

Database

Management

System

DATABASES: INTRODUCTION

I. Applied overview of Databases:-

The foundational reasoning "why we need database?" is to maintain a data store. So, our primary task is to build a mechanism in order to store data.

In order to get this, let us consider an example of a e-commerce website say "Amazon". Amazon stores humongous amount of data which includes information like:

→ Customer details, product details, purchases/returns, Inventory information, Shipment details, customer reviews, customer services etc.

This is just a small subset of data in Amazon and this data is being used by the amazon recommended system in order to predict the user choices and providing suggestions.

In order to design a database :

- (i) We need to understand the big picture :- we need to know how the pieces of data are related. For the previous example of amazon database, we should know how customers are related to other business logics like login, passwords, phone numbers etc. Note that we represent databases pictorially using E-R diagrams. E-R diagrams provides us this big picture.
- (ii) Tables to store the data :- Book-keeping and tables are the traditional ways to store the data on pre-computers. Since, tables are something which ^{is being} used from generations and is easy to understand. Therefore, database managers decided to store data in tabular form. Once we have the big picture, we need to find/know "how many tables will store this data in?" and "what should each table contain?" and that too occupying minimal storage in harddisk. The mathematical model which talks about how to store data in tables is called Relational model
- (iii) Storage of files on disk :- In order to store all data in computer we need to store it in disks. So, we need to store data in the disk such that we can access it optimally whenever required.
- (iv) Search/Retrieval fast :- The retrieving of data stored or searching of data should be fast and therefore, the database is stored in the form of B-Trees and B⁺-Trees.

(v) Multi-system access to your data:- There are situations when data is accessed by multiple systems / programs at a time simultaneously. So database manager should ensure that there is no flaw in the data during this concurrent access of data. This management system is termed as Transaction and Concurrency control system, whose job is to manage the concurrent access to the data.

NOTE:- language which is used to work on the database data is "SQL". SQL is a query language based on relational model and is used to add, delete, modify, search, retrieve data from the database.

II (E.2) Files vs DBMS :-

File System

→ Data is stored in disk in the format of comma separated files.

→ Querying:- Accessing any information from file system is a time consuming process.

DBMS

→ Data is being stored by database management system internally using indexing.

→ Querying:-
Database also stores info. in the form of file system but it is built based on special datastructure called Indexing
∴ Time required to access any information is complex.
(Objective of indexing is to make search

→ Redundancy:-

File systems contain/have redundant data and therefore consumes more disk space.

→ Redundancy:- +91 844-844-0102

It reduces redundancy, internally it avoids user to store the redundant data.

→ Consistency:-

It is difficult to have consistency of data in file system because interlinking of data is readily available in file system.

→ Consistency:-

Database management system is designed such that it takes care of the consistency of data internally because of interlinking of data.

→ Data independence:-

→ Data independence:-

DBMS provides a simple interface to query/search without thinking about how the data is stored. It hides the low level details from the engineers using DBMS.

→ Security and access control:-

→ Security and Access control:-

A DBMS can have multiple users, and each user have their access permissions, which were managed by DBMS internally.

→ Abstraction and ease of data access:-

→ Abstraction and ease of data access:-

It is the most important feature of DBMS. We need not care about how database works internally. We only ask questions using SQL.

Phone: +91 844-844-0102
 So, it abstracts away all the complexity of DS, algo, file storage.

NOTE: A filesystem can do whatever DB database does but for that we need to write 1000's of lines of code. In place of that we can simply install DB management system which does everything what file system does, plus huge other advantages.

(1.3) Tables and Keys:-

Terminology:

- 1) Relation:- Relation is nothing but a table. Table is a set of rows and columns.
 for example

CUSTOMER TABLE:

Column 1/ field 1/ attribute	CUST_ID	CUST_NAME	CUST_ADDR	CUST_PIN
1	abc	202,...	940861	→ row 1 / tuple 1
2	def	202,....	098910	→ row 2 / tuple 2
3	def	88, ne...	204102	⋮
4	xyz	civil lines..	509601	→ row 4

Fig: 1

The given customer table have 4 rows and 4 columns.

- 2) Attributes - Attributes are nothing but columns in relational database. These are sometimes also called as fields.

3) Tuples:- Rows are called as tuples in relational database. These are also known as records. In a relation there should not be two same tuples (tuple uniqueness constraint)

4) Instance:- The example given for customer table is an instance. That is, the whole table with some data filled in is called instance. "It is a set of tuples/rows along with the table structure."

5) Key:- It is the minimum number/set of attributes/columns to uniquely identify a row/tuple.

Consider the customer table, CUST_ID can uniquely identify a row. Therefore, CUST_ID can become the key of the relation.

- If we consider (CUST_ID, CUST_NAME) \rightarrow It can uniquely identify a row, but it is not a key because it is not minimum set of attributes. So, key has to be a minimum set of attribute to uniquely identify a row.

- If we consider (CUST_NAME, CUST_ADDR) \rightarrow It can uniquely identify a row and it is a key because neither CUST_NAME nor CUST_ADDR can uniquely identify a row separately.

6) Simple Key:- Key with only one attribute / column.

Ex- CUST_ID is the simple key for above relation.

7) Compound Key:- Key with multiple attributes / columns.

Ex- (CUST_NAME, CUST_ADDR) is a compound key.

8) Candidate key :- Set of all unique keys is called as candidate keys.

Ex - { CUST_ID, {CUST_NAME, CUST_ADDR} } are key candidates.

↓ Key 1 ↓ Key 2

Both Key 1 and Key 2 can uniquely identify a row/tuple.

9) Primary key :- One of the candidate key that the DB-designer chooses to maintain uniqueness is called primary key.

Ex - Either CUST_ID or {CUST_NAME, CUST_ADDR} could be chosen as primary key.

But there are certain rules while choosing a primary key.

If CUST_ID is chosen as primary key than

- It cannot be NULL for any tuple/row
- There should be at most one primary key per table
- Should be unique for each row/tuple.

These constraints / conditions are called as Entity integrity. Integrity constraints are properties to be satisfied while inserting, deleting or modifying the data in a relation/table.

NOTE - In a table/relation once a primary key has been chosen, it cannot have NULL values but, the remaining candidate key can have NULL values.

10) Alternate / secondary key: These are candidate keys that are not primary keys. That is they are unique to each row in but the DB designer do not choose it as primary, therefore, they can have NULL values now.

11) Relational Schema: Table / relation structure + Integrity constraints

12) Super key: Candidate key \cup other attributes.

Example $\{ \{CUST_ID\}, \{CUST_NAME\} \} \rightarrow$ it is a super key

$\{ \{CUST_NAME, CUST_ADDR\}, \{CUST_PIN\} \} \rightarrow$ it is also a Super key

13) Foreign key: It is an attribute or group of attributes in a relational database that provides a link between data in two tables.

Consider the below example :-

CUSTOMER RELATION

CUST_id	CUST_name	CUST_Addr	CUST_Pin
1	abc	202...	940861
2	def	202...	898910
3			
4			

(Referenced)

PURCHASES RELATION

Trans_id	CUST_id	Pg_id	Time
12345	1	12	5:00AM
45687	2	11	11:30
25678	5	14	18:00PM

foreign key

X
NOT VALID

Fig:2

(Referencing)

Every cust_id in PURCHASES RELATION has to be present in the CUSTOMER RELATION. The tuple (25678, 5, 14, 18:00) is not valid because CUST_id = 5 is not there in "CUSTOMER".

Hence, CUST_id is the foreign key, as it is a sort of pointer

Where each value of cust_id in PURCHASES Points to a value of CUST_id in CUSTOMER Relation, logically.

NOTE:- A foreign key may or may not be the primary key in referencing relation but it should be a primary key in referred relation, and hence, it removes redundancy and avoids inconsistencies.

- 14) Self-referential foreign keys:- A foreign key is said to be self referential, if an attribute or a group of attributes in one table that uniquely identifies a row of the same table. Consider the following example -

EMPLOYEE Relation:-

refers to			
Emp_id	Emp-name	Manager_id	DOJ
1	abc	3	06-06-96
2	:	:	:
3	:	:	:

Fig:3

Here, Emp_id is the primary key while Manager_id is a foreign key because it refers to the Emp_id (A manager is also an employee of the company). Therefore, Manager_id is a self-referential foreign key.

IV APPLIED
ROOTS

Integrity Constraints :-

1) Entity Constraint: It is as follows:

- ⇒ Every relation must have a primary key
- ⇒ Every primary key has to be unique
- ⇒ A primary key should not be NULL.

These 3 together are called entity constraints.

2) Domain Integrity Constraints: It gives/tells the column's/attribute's valid values.

Ex- For the Customer relation in fig: 1, Range of CUST_ID is $\rightarrow \{1, 2, \dots, 10,000\}$ and Domain is NUMERIC or INTEGER value.

3) User-defined Integrity Constraints: These constraints are business specific.

Consider the below example of Amazon's PURCHASES relation.

Trans_id	P_product	C_id	Time	# items
1206---7	Nike shoes	---	---	4

Phone: +91 844-844-0102

Amazon has set constraint over number of items (#items). You can buy of a product (say Nike shoes). It doesn't allow you to buy more than 4 pair of Nike shoes at a time. Such type of constraints are called user-defined Integrity constraint. which are business specific.

4) Referential Integrity constraint:-

As we can see in fig:2, CUSTID in PURCHASES is the foreign key which is referring to CUST-ID of CUSTOMER relation (which is the referenced relation). Referential integrity constraints corresponds to the foreign key concept and hence provides several constraints for inserting, deleting and modifying content in the referenced relation, CUSTOMER relation (in our set example) or referencing relation.

Changes to Referenced Relation:-

(i) On Inserting:- On inserting a new row in a referenced relation, no changes will be made in the referencing relation.

(ii) On Deleting:- On deleting a row from the referenced relation, changes will be made in the referencing relation based on the following set conditions.

(a) On delete no action:- On deleting 1st row from referenced CUSTOMER, no changes are made in PURCHASES. Though referential integrity is broken but we ignore that.

(b) On delete Cascade:- The moment we delete 1st row from CUSTOMER, on delete Cascade will make the cascaded changes in PURCHASES. that is , it will delete all row from PURCHASES referring to the deleted CUST_ID.

In our example on deleting tuple of CUST_ID=1 from CUSTOMER's will ^{also} delete 1st row from the PURCHASES due to cascade delete .

NOTE: It is a dangerous option to go for, therefore we avoid using it unless we know what we are doing. As this can result in loss of data.

On delete no action is also dangerous as it violates the referential Integrity.

(c) On delete Set NULL:- This is the safest option among all the three available options. It says, on deleting 1st row from CUSTOMER's , Set NULL in the CUST_ID CUST_ID is set to NULL for that particular referencing tuple.

NOTE:- There may be a condition where foreign key is a part of the primary key in referencing relation. In Such cases On delete Set NULL is not a feasible option to go for.

(iii) On Modifying / editing :- On modifying / editing a tuple is

in referenced relation, changes will be made in referencing relation based on following set condition.

(a) On modify no action:- If a row/tuple is modified in CUSTOMER, no action will be taken in the referencing relation. It is not widely used.

(b) On modify Cascade:- It is the safest option in case of modification of data in referenced relation. It will cascade the changes to the referencing relation.

(c) On modify set NULL:- It is not widely used.

It will set the referencing attribute in the particular tuple if any changes were made in referenced relation.

NOTE - On modify cascade is a time consuming operation because editing the changes in one relation will may also result editing the changes in other relation. Hence, it is time consuming.

NOTE - A relation cannot have two same tuples. every tuple should be unique (they should differ by atleast one attribute ie. primary key)

ex - R,

A	B
a ₁	b ₁
a ₂	b ₁

some { a₁ b₁ } X → This is not possible as (1) and (2) are same



Changes to the Referencing Relation:-

(i) On inserting :- On inserting any tuple/row in referencing relation PURCHASES, first it will check CUSTOMER's relation if the inserted CUST-ID exists in CUSTOMER or not. In Fig:2, the entry made in PURCHASES with CUST-ID=5 is not valid because 5 does not exists in CUSTOMER's and hence, it will throw an error.

Therefore, on inserting check for any violation.

(ii) On deleting :- On deleting a row from referencing relation will make/brings no change in the referenced relation

(iii) On modifying/editing : We need to check for any violations regarding foreign key constraint in referenced relation. If any, DBMS will raise an error.

Example - What all the rows getting deleted on deleting (2,4) from R
 Let R be a relation given as follows and DB manager set the constraint as ON DELETE CASCADE.

S.N.O.	A ₁	A ₂
1	2	4
2	3	4
3	4	3
4	5	2
5	7	2
6	9	5
7	6	4

⇒ Since, the constraint is ON DELETE CASCADE.

Therefore, on deleting (2,4)

- * Initially ~~both~~ 4th and 5th tuple will get deleted as they are referring to '2' so (5,2) and (7,2) will get deleted.

- * Deletion of (5,2) and (7,2) will lead to the multilevel cascading, so all tuples referring to 5 and 7 will also get deleted.

Therefore, total number of tuples deleted from R, on deleting (2,4) are 3, and they are :- (5,2), (7,2) and (9,5)

This is the reason why on delete cascade is a dangerous option as it leads to the multilevel cascading and hence, there is loss of data.

V (1.5) Solved Problems:-

(81) $R(A_1, A_2, A_3, \dots, A_n)$ (This represents a relation R having n attributes)

(a) Number of superkeys possible, if candidate keys are $\{A_1\}$?

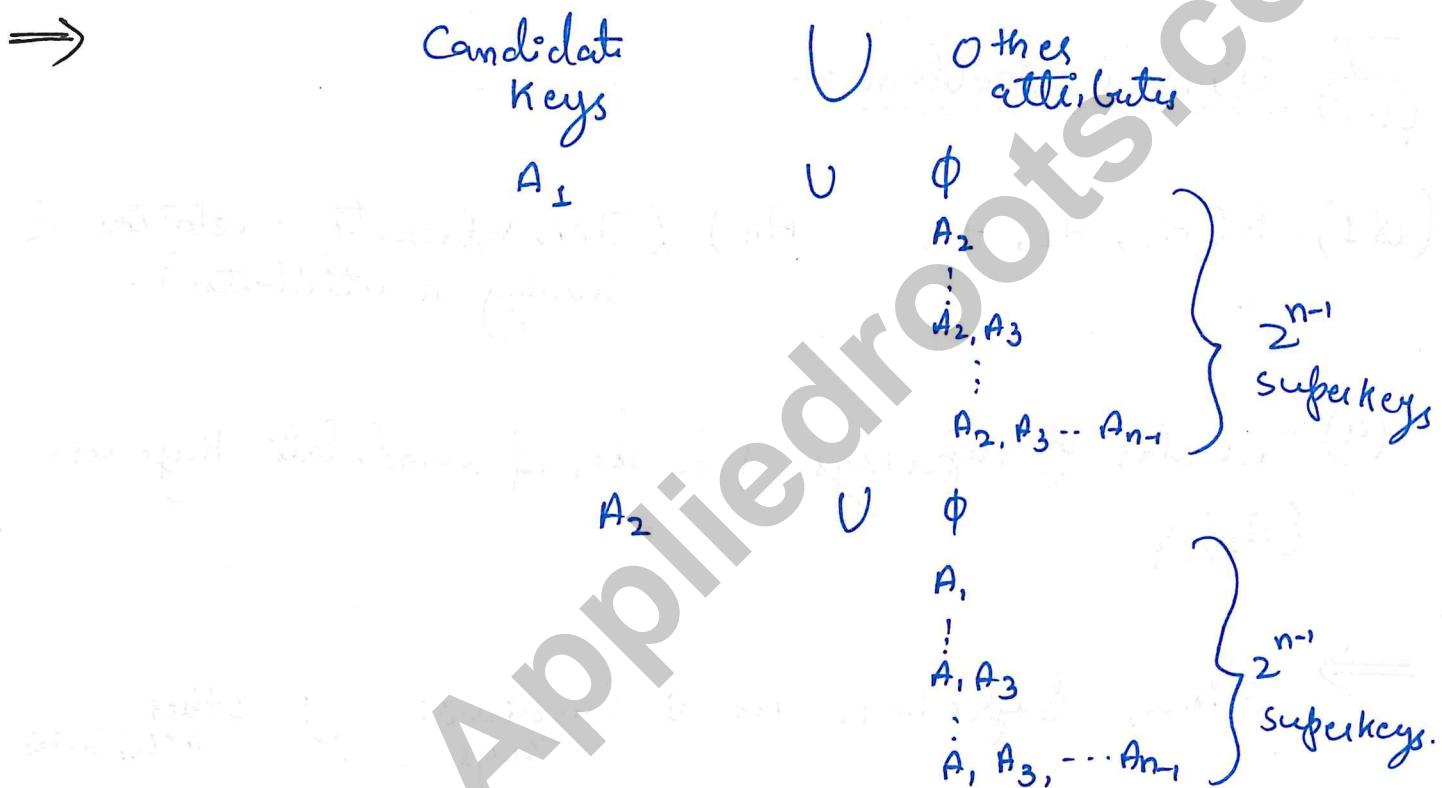
Since, Super keys are \Rightarrow Candidate Keys	\cup	Other attributes
A_1	\cup	\emptyset
it is candidate key given in question.	A_2	A_2
	A_3	A_3
\vdots	\vdots	\vdots
	$A_2 A_3$	$A_2 A_3$
\vdots	\vdots	\vdots
	$A_2, A_3 \dots A_n$	$A_2, A_3 \dots A_n$

So A_1 , can be merged with remaining $(n-1)$ items in relation R. So this is nothing but the concept of power set. for a set of n elements, no. of subset we can create using n elements are 2^n . Similarly, here for $(n-1)$ elements (bcz A_1 is already taken) we can have subset = 2^{n-1}

When each of these 2^{n-1} subsets get combined with A_1 it will become a super key.

Hence, # super keys given candidate key as $\{A_1\}$ are 2^{n-1}

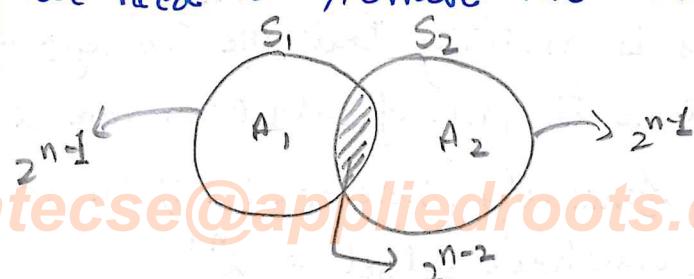
(b) For given relation R, what is the number of super keys if, candidate keys are $\{A_1, A_2\}$



As we can see, we have counted set $\{A_1, A_2\}$ twice as it appears in both A_1 and A_2 . Similarly there are many other subsets which we have counted twice.

(Note Set $\{A_1, A_2\} \Leftrightarrow \{A_2, A_1\}$)

So Acc. to the principle of inclusion and exclusion from set theory we need to remove the intersection part.



So, $|S, US_2| = |S_1| + |S_2| - |S, NS_2|$ Phone: +91 844-844 0102 (acc. to the inclusion-exclusion principle)

These are those superkeys which are counted twice and contains A_1, A_2 together.

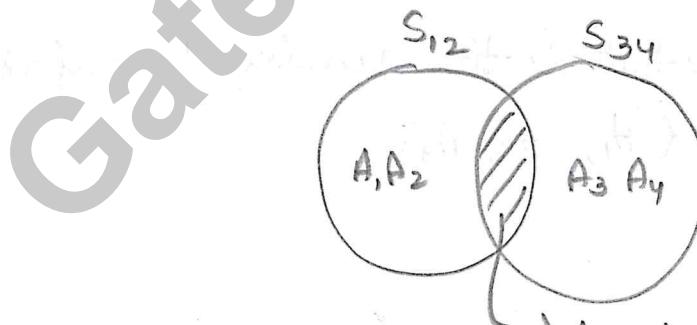
Therefore, # of superkeys = $2^{n-1} + 2^{n-1} - 2^{n-2}$

$$\begin{aligned} &= 2 * 2^{n-1} - 2^{n-2} \\ &= 2^n - 2^{n-2} \end{aligned}$$

(C) For given relation R, what is the number of superkeys if candidate keys are $\{\{A_1, A_2\}, \{A_3, A_4\}\}$

\Rightarrow With $\{A_1, A_2\}$ as candidate key } = 2^{n-2}
 Number of subsets will be (S_{12})

With $\{A_3, A_4\}$ as candidate key } = 2^{n-2}
 Number of subsets will be (S_{34})



A_1, A_2, A_3, A_4 are the elements in this intersection. So No. of subsets = 2^{n-4}

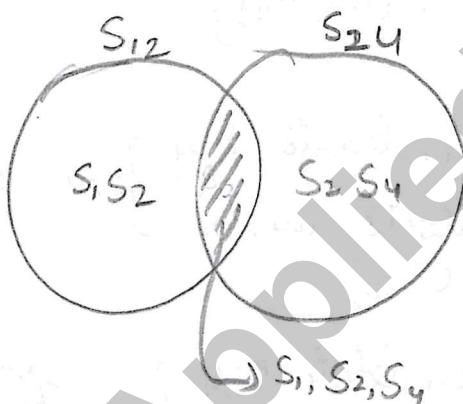
Therefore, no. of superkeys = $2^{n-2} + 2^{n-2} - 2^{n-4}$

(d) For the given Relation R, what is the number of superkeys if, candidate keys are $\{ \{A_1, A_2\}, \{A_2, A_4\} \}$

\Rightarrow For $\{A_1, A_2\}$ as candidate key, no. of subset will be $\{S_{12}\} = 2^{n-2}$

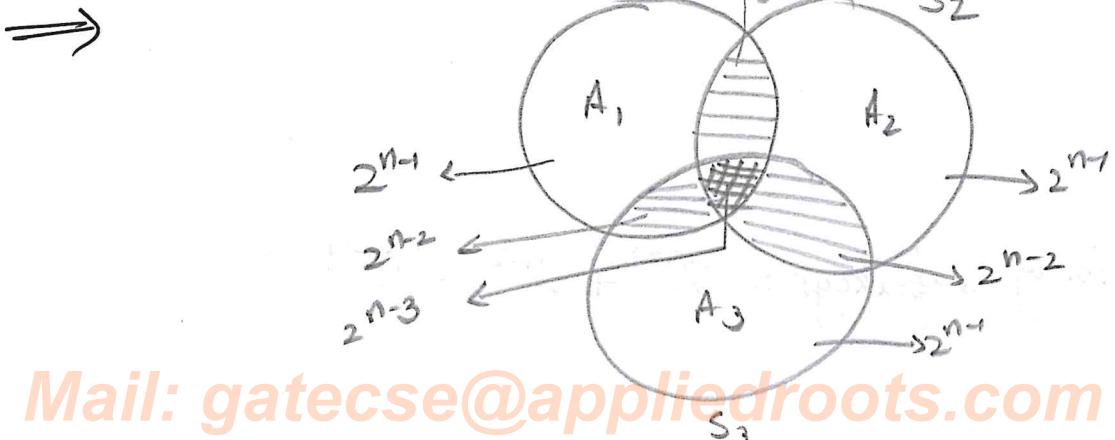
For $\{A_2, A_4\}$ as candidate key, no. of subset will be $\{S_{24}\} = 2^{n-2}$

For $\{A_1, A_2, A_4\}$ no. of subsets will be $= 2^{n-3}$



Therefore, no. of superkeys = $2^{n-2} + 2^{n-2} - 2^{n-3}$

(e) For given relation R, what is the number of superkeys if, candidate keys are $\{A_1, A_2, A_3\}$



According to the principle of inclusion and exclusion

$$|S_1 \cup S_2 \cup S_3| = |S_1| + |S_2| + |S_3| - |S_1 \cap S_2| - |S_1 \cap S_3| - |S_2 \cap S_3| + |S_1 \cap S_2 \cap S_3|$$

$$\begin{aligned} \text{Therefore, no. of Super keys} &= 2^{n-1} + 2^{n-1} + 2^{n-1} - 2^{n-2} - 2^{n-2} \\ &\quad - 2^{n-2} + 2^{n-3} \\ &= 3 \cdot 2^{n-1} - 3 \cdot 2^{n-2} + 2^{n-3} \end{aligned}$$

(Q2) Given a relation R with n attributes, how many superkeys can be present?

\Rightarrow Given, R (A₁, A₂, ..., A_n)

Since, no information is given regarding candidate keys.
Therefore any attribute could be the key

If, A₁ is key \rightarrow # superkeys = 2ⁿ⁻¹

If, A₂ is key \rightarrow # superkeys = 2ⁿ⁻¹

⋮

So, total 2ⁿ subsets are possible except NULL set({}) can be a super key. (\emptyset could not be the superkey)

So, # superkeys will be = 2ⁿ - 1 superkeys

↳ because {} cannot be a superkey.

ENTITY RELATIONSHIP (ER) MODELS & DIAGRAMS

I Introduction to ER Diagrams :-

(2.1)

Purpose / objective of ER diagram or model is to get a high level picture of the database. It gives the logical view of the database.

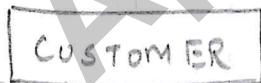
What is an Entity ?

Entity is an object that can be uniquely identified.

For ex:- In fig: 2, Customer and Purchases are two entities as these are the physical or logical concept.

In an E-R diagram, an entity is represented as a simple rectangular box.

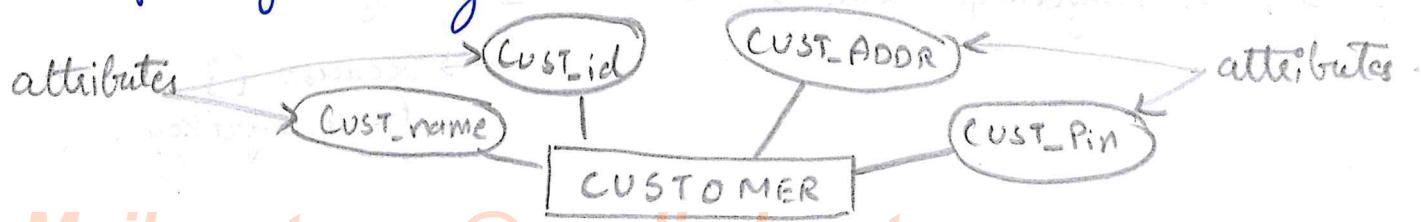
ex →



→ Attributes :- Attributes are properties / features of an entity. It is typically represented using ellipse.

for example -

On considering fig: 1, we can represent in the form of ER diagram as:



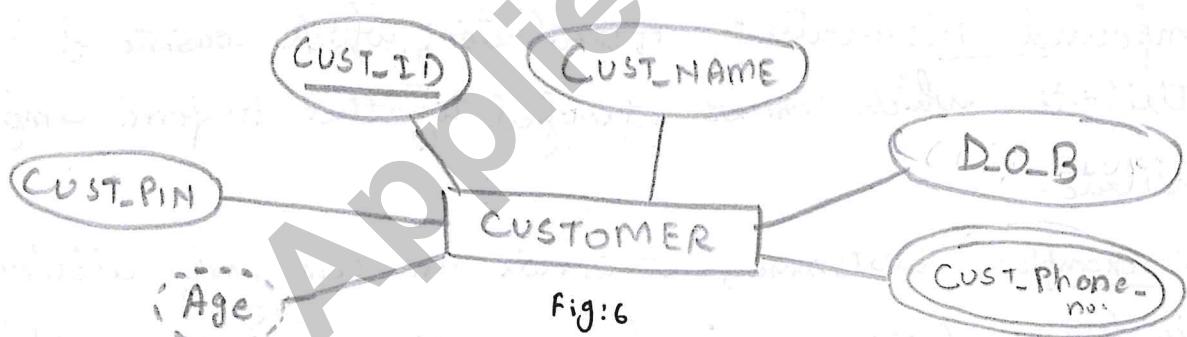
NOTE: A relationship can also have attributes, those attributes are called declarative attributes.

NOTE: The attributes of an entity are nothing but column names / attribute names in relational model.

An each row in the relation / table represents an entity. For the customer example in fig:1, each row represents a customer.

→ key Attributes:- Key attributes are the one which uniquely identifies a row. It is often represented as an ellipse and is underlined.

For example: In fig:6, CUST-ID is the key attribute.



→ Multi-valued Attributes:- Attributes which are having multiple values associated with them are multi-valued attributes. These are typically represented by two concurrent ellipses.

For example In fig:6, CUST_Phone_no. is a multi-valued attribute. A customer can have multiple contact numbers.

→ Derived Attributes :- An attribute whose value is derived / calculated for other attributes are called derived attributes. These are represented by dotted ellipse.

For example: In fig:6, age is an derived attribute, which is calculated / derived from D.O.B (date of birth of a customer)

~~Note~~ NOTE:- Derived attributes are not typically stored within the database, physically ; instead it can be derived by using an algorithm.

→ Compound Attributes :- Attributes, which consists of other attribute and can be grouped together to form compound attribute.

For example: CUST-name, CUST-addr are compound attributes in fig:7. CUST-name is formed using first-name, middle-name, and last-name. So all of these constitutes to form the CUST-name.

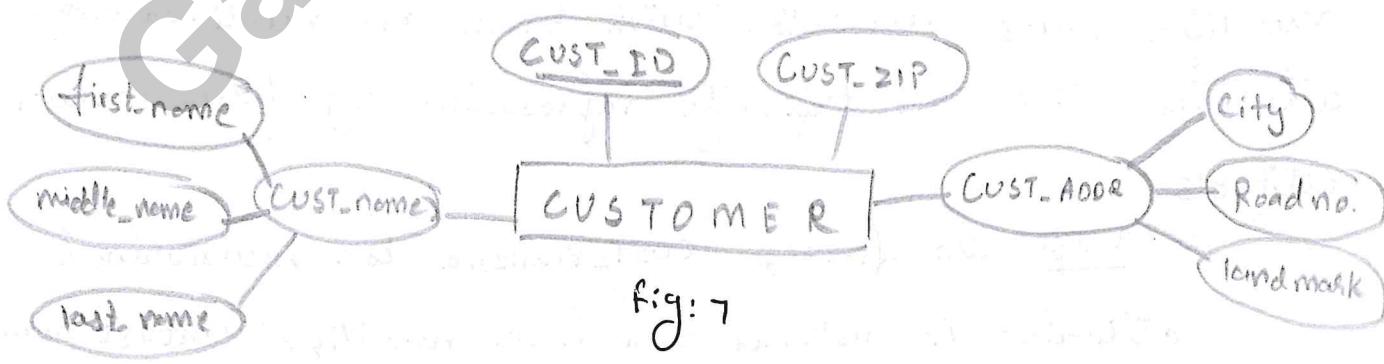
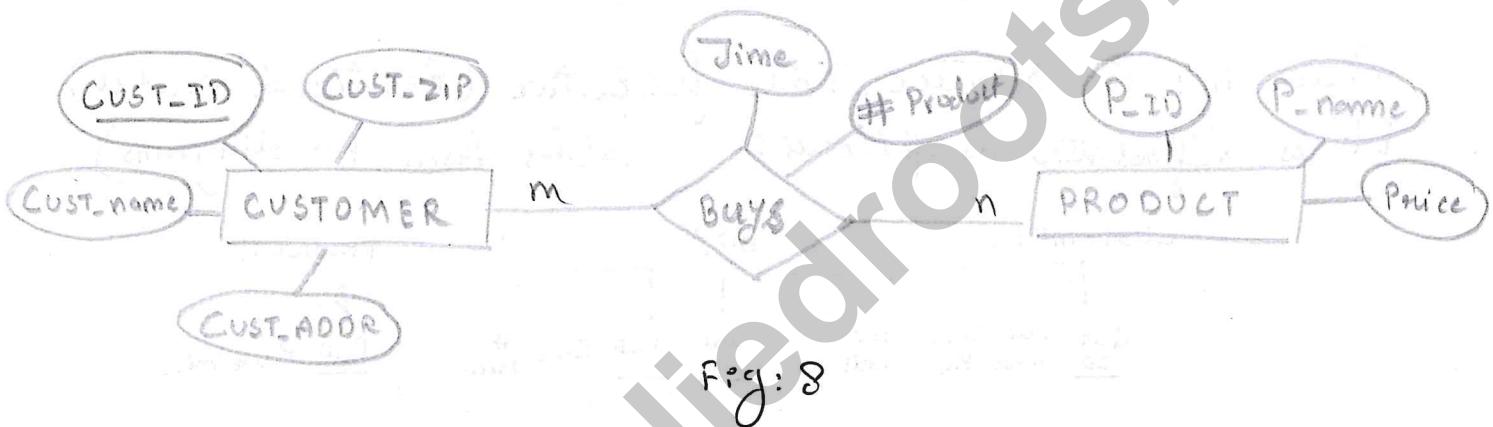


fig:7

Similarly, CUST-ADDR is a compound attribute as it is formed of City, landmark and Road no.

Relationships :- A relationship is a situation (in context of databases) which exists in between multiple entities. It is typically represented using a rombus / diamond.

for example:- In fig:8, "buys" is a relationship between Customer and products. This could be read as "A customer buys a product/s". Therefore, relationship is a way to connect multiple entities.



When a customer buys a product, there is some time stamp associated with it. So, this attribute "time" will be the attribute of "Buys" and not of CUSTOMER or PRODUCT because for us to have the time, we need to have both CUSTOMER and a PRODUCT. Similarly, #Product will be the attribute of relationship "Buys".

According to set theoretic perspective, (NOTE: all of the database concept can be understood from the concept of sets and relation) it will be represented as follows:

Every dot
represents
a customer

CUSTOMER

PRODUCTS

Phone: +91 844-844-0102

Every dot
represents
a product

Fig:9

Acc. to the relation depicted above, many customers can buy a product, many products could be bought by a single customer or many products could be bought by many customers.

According to relation/Table perspective (Relational model) it is represented as follows: (Tables from E-R diagrams)

CUSTOMER:				BUYS:				PRODUCT:		
CUST ID	CUST name	CUST Pin	CUST add	CUST ID	P_ID	Time	# Prod.	P_ID	P_name	Price

Fig:10

Since, the given relation is many-to-many in between CUSTOMER and PRODUCTS, therefore, we need a separate table/ relation for "buys". Otherwise it would be possible to store all information using two relations.

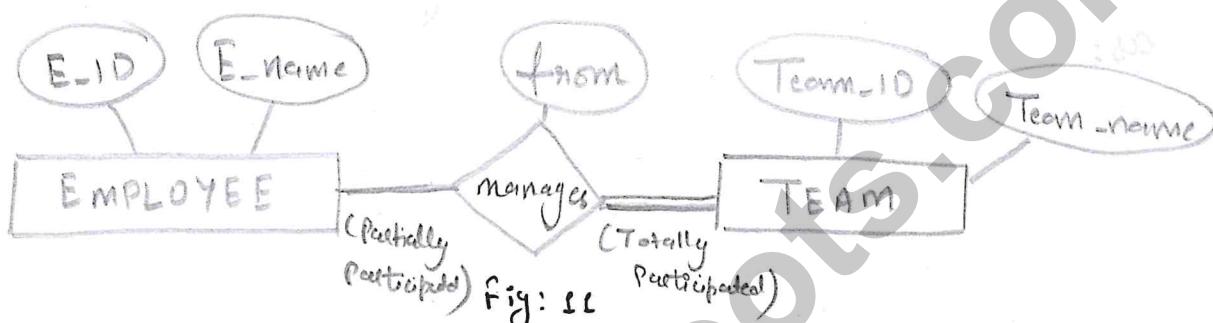
(This concept is further explained in the detail)

Note, In BUYS relation, CUST ID and P_ID are the foreign key to CUSTOMER and PRODUCTS.

Mail: gatecse@appliedroots.com

→ Total and partial participation:-

This participation is about entities participating in a relationship. Consider the below E-R diagram of 'Employee' and a 'Team' related to each other using relationship "manages".



For the given E-R diagram in Fig:11, the diagram could be explained as "Every team needs to have a manager and a Employee/ manager may manages a team from some give date".

So, here total participation could be understood as - For every 'Team' there should be a manager OR every member of Entity 'Team' should participate in the relationship with Entity 'Employee'. Hence, 'Team' is said to be totally participated with the "manages" relationship.

Note: Total participation is represented using two parallel lines (as we can see in figure fig:11)

While, entity Employee does not participated totally with relationship manages because every Employee is need not suppose to be a manager.

Since, there are employee who do not manage any team, therefore, entity employee is having partial participation with the relationship "manages". It is represented using a single line (it is set default).

Based on set theoretic concept it could be represented as:

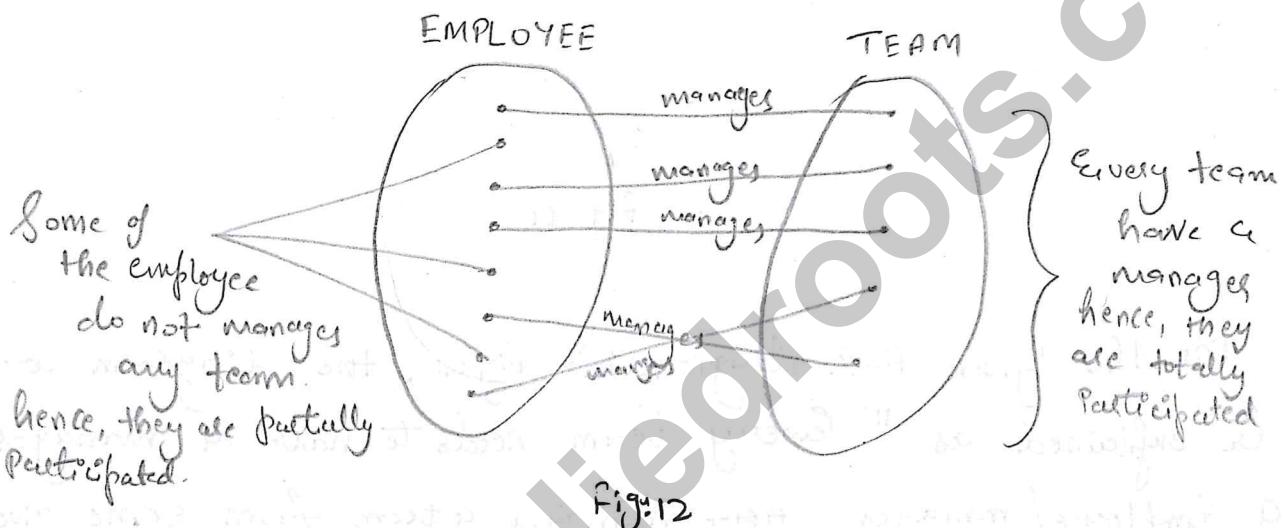


Fig:12

Note: We can think of total and partial participation as a form of user-defined Integrity constraints. (As it depends as per business requirement)

II

Cardinality of Relationships and Constructing

(2.2)

Minimal Tables / Relations :-

The 4 possible cardinalities of Relationship are:-

- One-to-One
- One-to-many
- Many-to-One
- Many-to-many

→ One-to-One Relationship :- (1:1)



Fig: 13

One-to-One relationship is represented using an 'arrow': (\rightarrow, \leftarrow)

In Fig:13, as we can see both customer and driver's license is having an arrow which means "The entity "customer" is participated in the relation "has a" in a one-to-one fashion" and it could be read as:

"Each customer has exactly one driver's license and every driver's license is associated with exactly one customer"

Fig:14, depicts the set theoretic perspective of E-R diagram given in Fig:13.

There are some customers who has not uploaded their driver's license

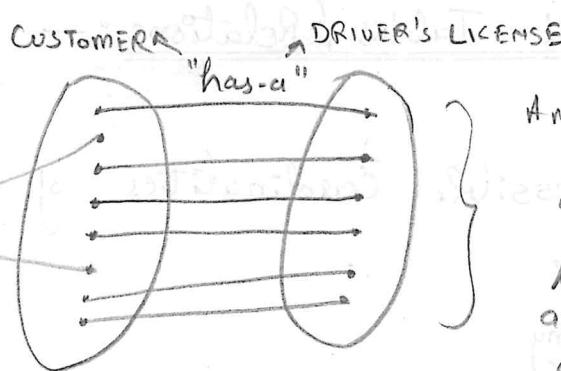


Fig:14

As, we can see in Fig:14, DL is totally participated with "has-a" while customers are partially participated with relationship "has-a".

→ One-to-many Relationship :- (1:m)

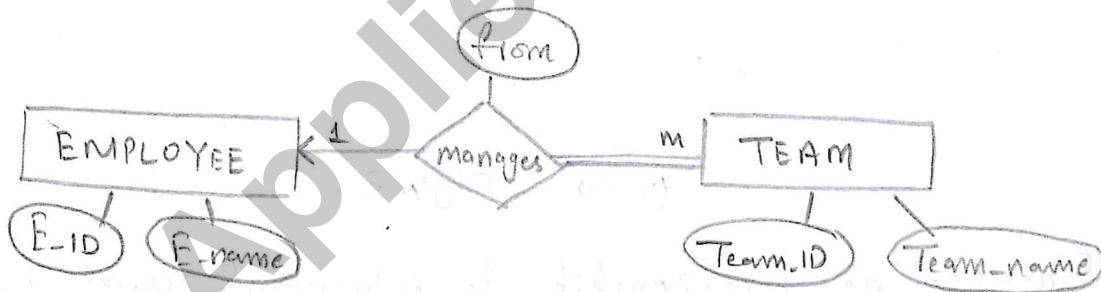
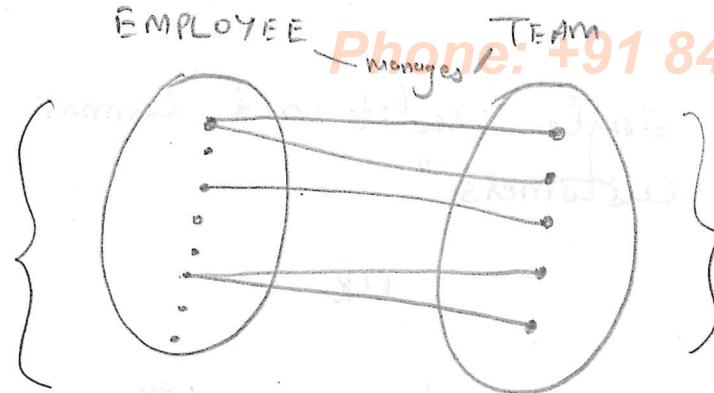


Fig:15

Single or double line by default depicts the many part of the relationship. This relation is one-to-many because "there are some employees who manages more than one team while there is only one manager to manage a team".

Fig:16 depicts the set theoretic perspective of E-R diagram shown in Fig:15.

A employee can manage more than 1 team and there are some employees who do not manage any team.



Every team is managed by exactly one employee.

Fig:16

NOTE:- Any of these 4 cardinalities of relationship can have both side factual, both side total, one side factual and one side total participation. They can have any permutation or combination of possibilities of a E-R diagram.

for ex- One-to-one can have any of the 3 possible condition of participation in an E-R diagram.

→ Many-to-one Relationship :- (m:1)

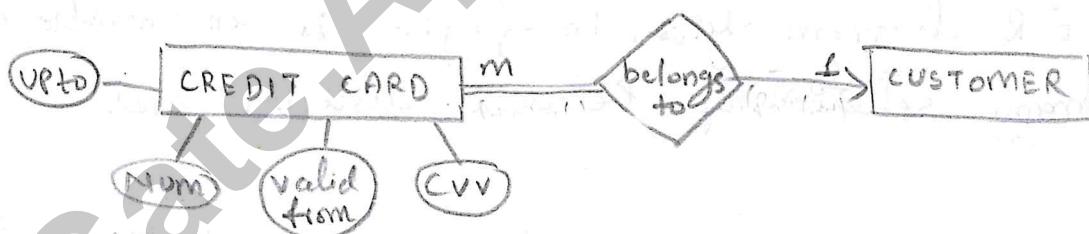


Fig:17

It is related to one-to-many relationship (the difference is due to the perspective).

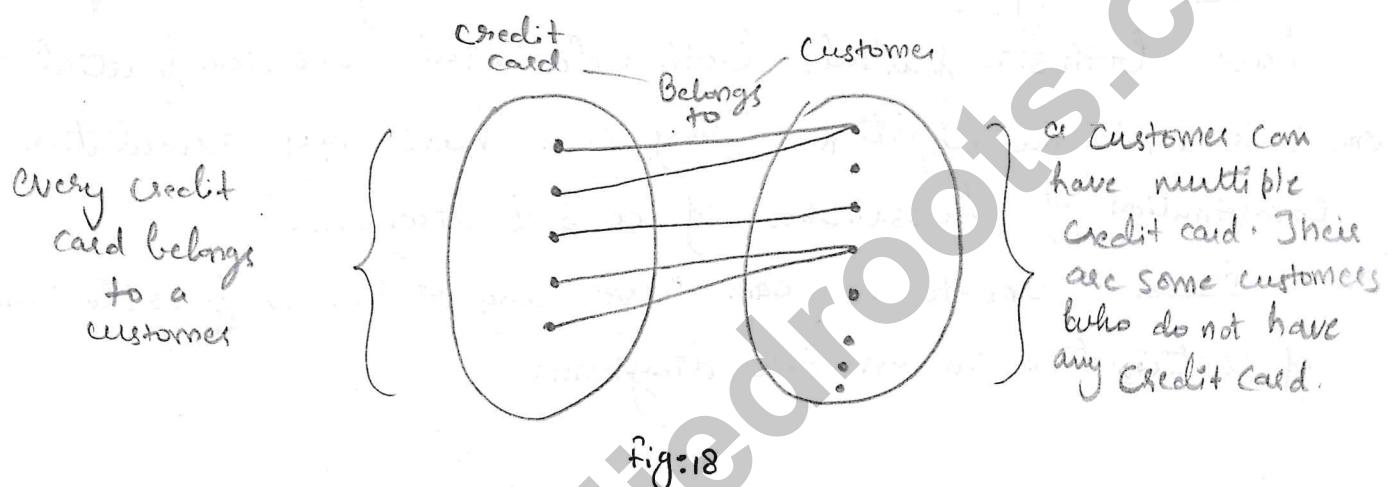
" Any ^(multiple) number of credit cards ^{can} belongs to a customer (and there might be some customer who do not have any credit card but every credit card we have in our system has to be associated with an employee (total Participation))

APPLIED ROOTS
while a single credit card cannot belongs to the
multiple customers"

OR

"A customer can have multiple credit cards but
a credit card cannot have multiple customers"

Fig:18, depicts its set theoretic perspective :-



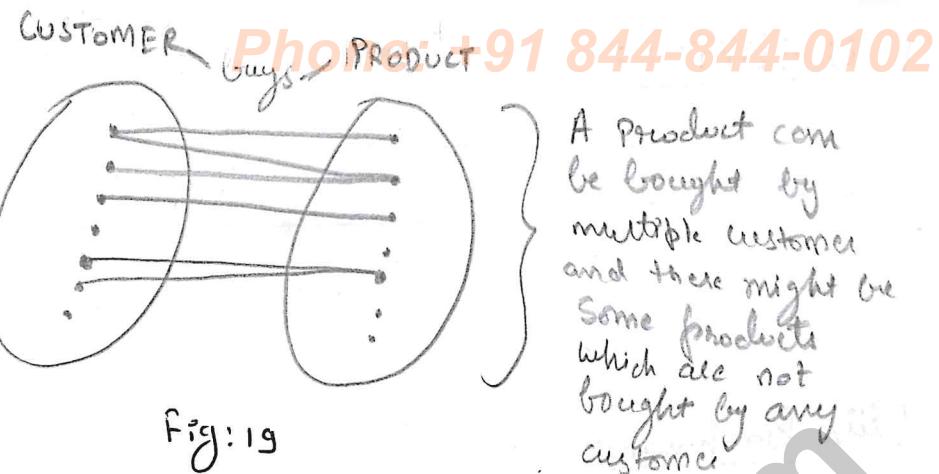
→ ↵ Many-to-Many Relationship :~ (m:n)

The E-R diagram given in fig:8 is an example of many-to-many relationship between Customer and Product.

"A customer can buy many/multiple products and a product could be bought by many/multiple customers".

Fig:19 depicts its set theoretic perspective : -

A customer can buy multiple products and there are some customers who didn't buy any product.



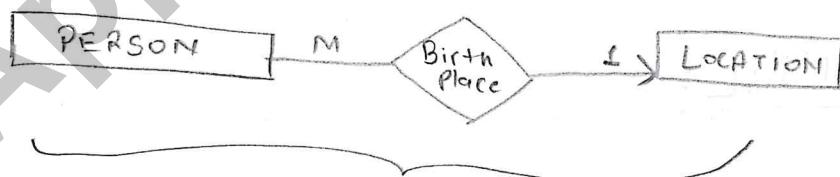
A product can be bought by multiple customers and there might be some products which are not bought by any customer.

→ Alternative Conventions :-

What we have used till now to represent the cardinality of a relationship is :-

- To depict 'one'
- To depict 'many'

for example: For many-to-one :-



This is one of the way to represent cardinality.

Other popular alternative conventions are:-

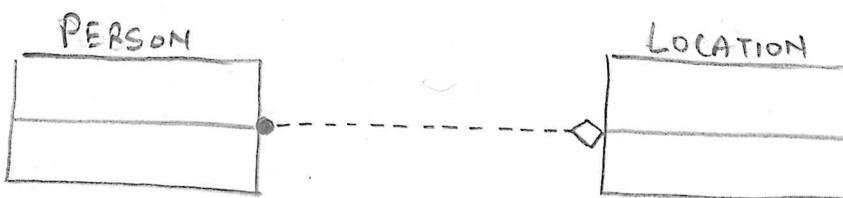
(i) Chen:



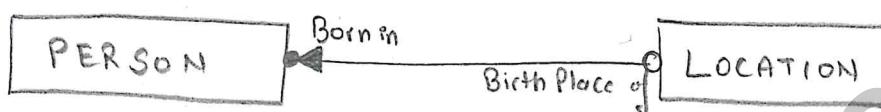
(ii) IDEFIX:

Phone: +91 844-844-0102

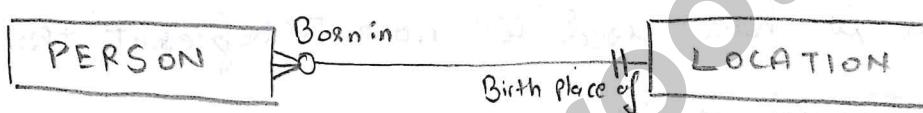
APPLIED
ROOTS



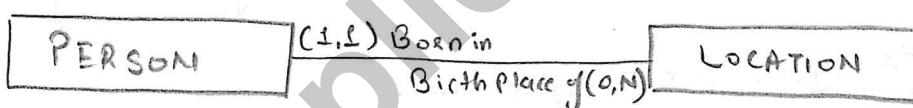
(iii) Bachman:



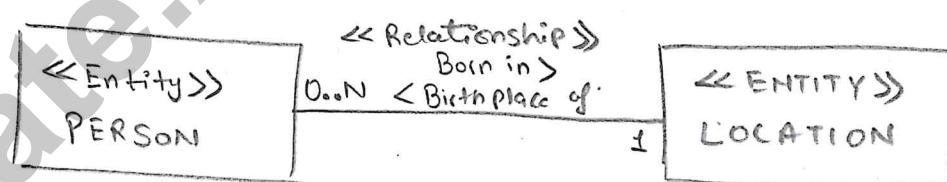
(iv) Martin / IE / Crow's Foot:



(v) Min-Max / ISO:



(vi) UML:



NOTE:- All of the above alternative Conventions are the example for many-to-one relationship.

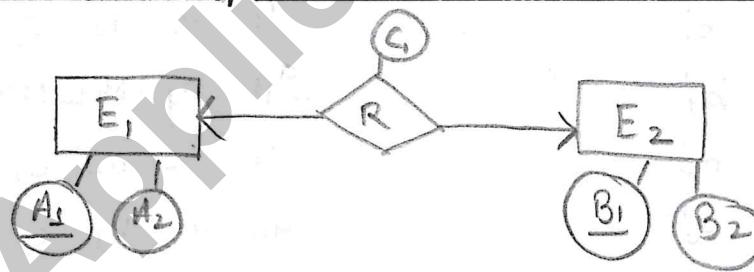
(8) What is the minimum # relations for every possible combination of cardinality and participation of between ~~the~~ entries? **Phone: +91 844-844-0103**

→ Constructing tables from E-R Diagrams :-

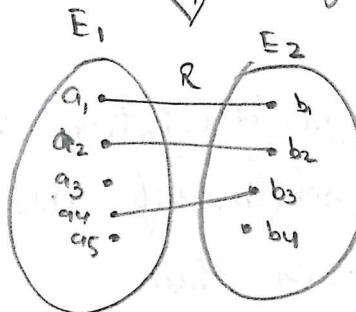
Given an E-R diagram we need to look for cardinality, participation, keys etc in order to construct tables from that. Along with this there is some set of rules defined for constructing a table / relation. These are as follows:

- [1] We need to minimize the number of table / relations.
- [2] Every relation / table should have a primary key.
- [3] Each cell in table / relation must have only one value.
(Cells can have NULL value except the primary key cells)

(i) One-to-one relationship + Partial-Partial participation :-



↓ Set theoretic diagram



Check for f Relation:

A ₁	A ₂	B ₁	B ₂	C ₁
a ₁	a _{1'}	b ₁		
a ₂	a _{2'}	b ₂		
a ₃	a _{3'}	NULL		
a ₄	a _{4'}			
a ₅	a _{5'}			

⇒ Check for 1 relation is possible or not: Phone: +91 844-844-0102

A ₁	A ₂	B ₁	B ₂	C ₁
a ₁	a' ₁	b ₁	b' ₁	c ₁
a ₃	a' ₃	NULL	NULL	NULL
NULL	NULL	b' ₂	b' ₄	NULL

cannot become key

As we can see, all of them contain NULL values hence neither A₁ nor B₁ could be the primary key. Hence, it is not possible to set them into 1 relation. (It is Violating the 2nd rule [2])

⇒ Check for 2 relation is possible or not:

Table 1 -

B ₁	B ₂
b ₁	b' ₁
b ₄	b' ₄
b ₂	b' ₂
b ₃	b' ₃

Key values
are not NULL

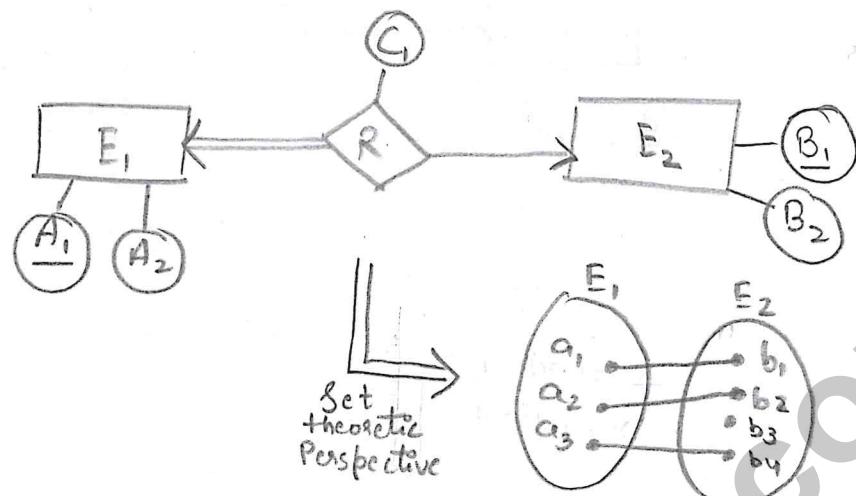
refer Table 2 -

A ₁	A ₂	B ₁	C ₁
a ₁	a' ₁	b ₁	c ₁
a ₃	a' ₃	NULL	NULL
a ₂	a' ₂	b ₂	c ₂
a ₄	a' ₄	b ₃	c ₃
a ₅	a' ₅	NULL	NULL

Key values
are not NULL

Hence, minimum no. of relations required for one-to-one, both side factual relationship are two. And maximum number of relation is obviously three.

(ii) One-to-One relationship + at least one side total Participation



⇒ Check for 1 relation is possible or not:

A ₁	A ₂	B ₁	B ₂	C ₁
a ₁	a ₁ '	b ₁	b ₁ '	c ₁
a ₂	a ₂ '	b ₂	b ₂ '	c ₂
a ₃	a ₃ '	b ₃	b ₃ '	c ₄
NULL	NULL	b ₃	b ₃ '	c ₃

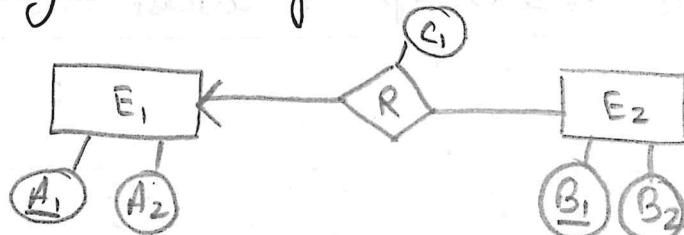
↓

cannot become key B₂ can become key for this table

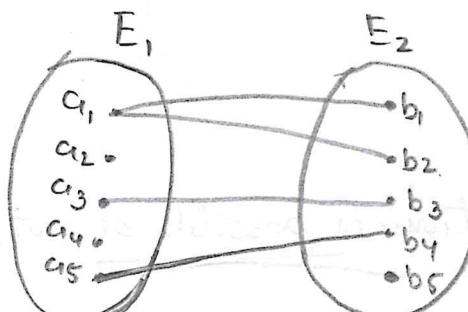
As, we can see there is no NULL value for attribute B₂, therefore, it B₂ can become the primary of this table. Hence, minimum no. of relation for given condition is 1.

Note:- If the case given is one-to-one + Both side total then definitely one relation is required to represent it, and the key for that relation could be either A₁ or B₂ as none of the attributes can have NULL values.

iii) One-to-Many Relationship + Partial-Partial Participation - 0102



↓ Set theoretic perspective



⇒ Check for 1 relation is possible or not:

A ₁	A ₂	B ₁	B ₂	C ₁
a ₁	a ₁ '	b ₁	b ₁ '	c ₁
a ₁	a ₁ '	b ₂	b ₂ '	c ₂
a ₂	a ₂ '	NULL	NULL	NULL
a ₃	a ₃ '	b ₃	b ₃ '	c ₃
a ₄	a ₄ '	NULL	NULL	NULL
a ₅	a ₅ '	b ₄	b ₄ '	c ₄
NULL	NULL	b ₅	b ₅ '	NULL

None of these could become primary key.

Since, all the attributes has violated the [2] condition hence, 1 relation is not possible.

Phone: +91 844-844-0102

⇒ Check for 2 relations are possible or not :-

APPLIED ROOTS

Table 1:	<u>A₁</u>	<u>A₂</u>
	a ₁	a ₁ '
	a ₂	a ₂ '
	a ₃	a ₃ '
	a ₄	a ₄ '
	a ₅	a ₅ '

refers to ↓

Table 2:	<u>B₁</u>	<u>B₂</u>	<u>C₁</u>	<u>A₁</u>
	b ₁	b ₁ '	c ₁	a ₁
	b ₂	b ₂ '	c ₂	a ₁
	b ₃	b ₃ '	c ₃	a ₃
	b ₄	b ₄ '	c ₄	a ₅
	b ₅	b ₅ '	NULL	NULL

Two relations will work as there are no NULL values for A₁ and B₂ in relation 1 and 2 respectively.

Hence, 2 relations are required in minimum for the given condition.

Ques

NOTE: There are many ways to split the table but all splits will not work. What would be the result if the above relations are not splitted like this. Let us consider the follow two relations:

Table 1:

	<u>A₁</u>	<u>A₂</u>	<u>B₁</u>	<u>C₁</u>
	a ₁	a ₁ '	{b ₁ , b ₂ }	c ₁
	a ₂	a ₂ '	NULL	NULL
	a ₃	a ₃ '	b ₃	c ₃
	a ₄	a ₄ '	NULL	NULL
	a ₅	a ₅ '	b ₄	c ₄

Table 2:

	<u>B₁</u>	<u>B₂</u>
	b ₁	b ₁ '
	b ₂	b ₂ '
	b ₃	b ₃ '
	b ₄	b ₄ '
	b ₅	b ₅ '

↓
this is
having
2 values

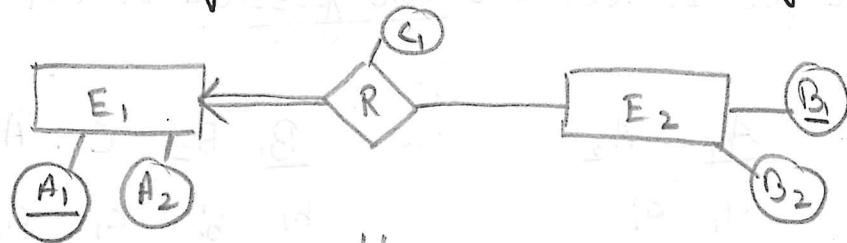
It violates the [3] rule and therefore such split of relations is not possible.

Mail: garcse@appliedroots.com

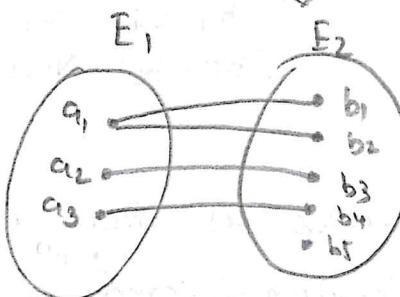
(iv) One-to-many relationship + One side total participation 14-0102

**APPLIED
Roots**

Case 4:
One-to-many
+
total participation
"one" side.



↓ Set theoretic
Perspective



⇒ Check for 1 relation is possible or not:

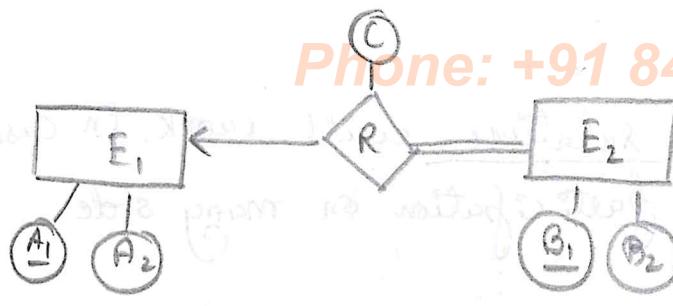
A ₁	A ₂	B ₁	B ₂	C ₁
a ₁	a ₁ '	b ₁	b ₁ '	c ₁
a ₁	a ₁ '	b ₂	b ₂ '	c ₂
a ₂	a ₂ '	b ₃	b ₃ '	c ₃
a ₃	a ₃ '	b ₄	b ₄ '	b ₄
NULL	NULL	b ₅	b ₅ '	NULL

↓
This could be the
Primary key

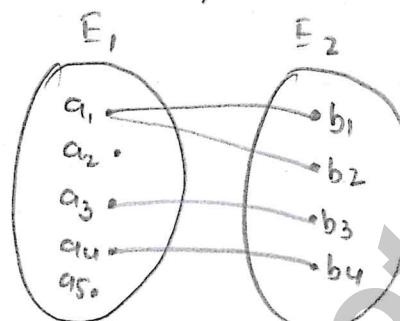
Hence, 1 relation is sufficient for one-to-many and total participation is on the "one" side.

Case 2:

One-to-many
+
total participation
on many sides



↓
Set diagram



→ check for 1 relation is possible or not:-

A ₁	A ₂	B ₁	B ₂	C ₁
a ₁	a' ₁	b ₁	b' ₁	c ₁
a ₁	a' ₁	b ₂	b' ₂	c ₂
a ₂	a' ₂	NULL	NULL	NULL
a ₃	a' ₃	b ₃	b' ₃	c ₃
a ₅	a' ₅	NULL	NULL	NULL

Neither A₁ nor B₁ can become the primary key.

Hence, 1 relation is not possible

→ check for 2 relation is possible or not:-

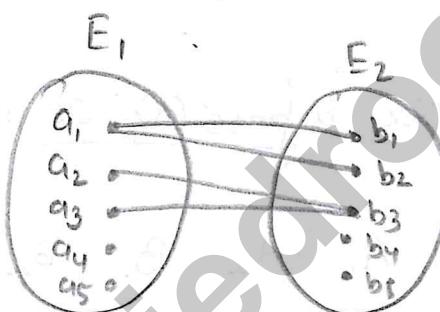
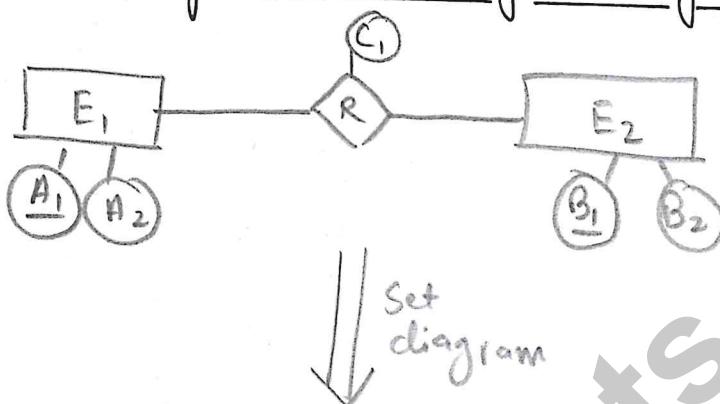
A ₁	A ₂	B ₁	B ₂	C ₁	A ₁
a ₁	a' ₁	b ₁	b' ₁	c ₁	a ₁
a ₂	a' ₂	b ₂	b' ₂	c ₂	a ₁
a ₃	a' ₃	b ₃	b' ₃	c ₃	a ₃
a ₄	a' ₄	b ₄	b' ₄	c ₄	a ₄
a ₅	a' ₅				

Hence, 2 relations will work in case of one-to-many and total participation on many side.

(v) Many-to-Many Relationship + Both side partial participation:-

Case 1:-

Many-many
+
Both partial.



We can clearly see that neither 1 relation nor 2 relation will work in this case as each a_i is related to 2 values of E_2 and b_j is related to two values of E_1 . So, in which ever way we break this down in 2 tables it will result in some invalid tuples. Hence, 2 relations will not work.

In this case we need 3 relations minimum, two relations to store E_1 and E_2 while third relation is to store relation R (along with its relationship attributes and primary keys of E_1 and E_2)

1)

<u>A₁</u>	<u>A₂</u>
----------------------	----------------------

2)

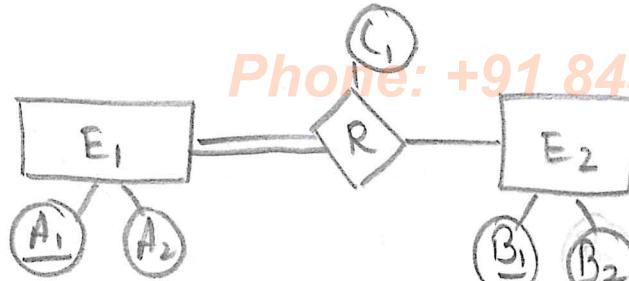
<u>C₁</u>	<u>A₁</u>	<u>B₁</u>

3)

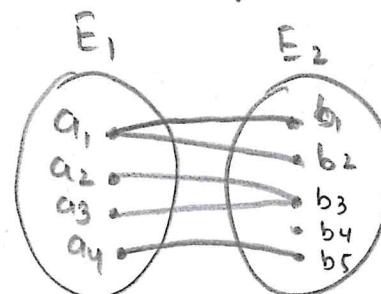
<u>B₁</u>	<u>B₂</u>
----------------------	----------------------

Case 3:

Many-many
+
one side
total
Participation



↓
Set diagram



⇒ Check for 2 relation is possible or not :-

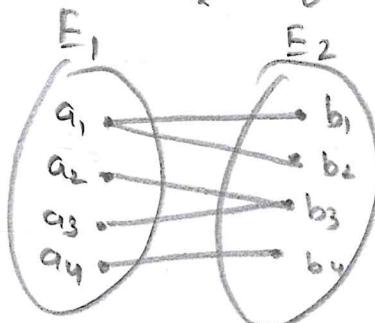
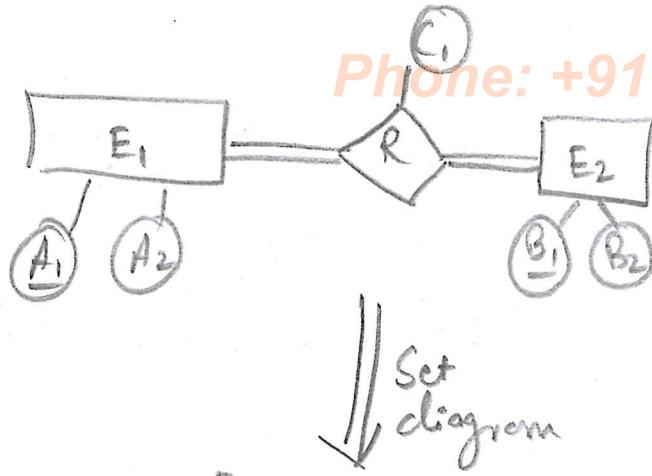
<u>A₁</u>	<u>A₂</u>	<u>B₁</u>	<u>C₁</u>		<u>B₁</u>	<u>B₂</u>
a ₁	a' ₁	{b ₁ , b ₂ }	c ₁	X	b ₁ b ₂ b ₃ b ₄ b ₅	b' ₁ b' ₂ b' ₃ b' ₄ b' ₅

OR

<u>A₁</u>	<u>A₂</u>		<u>B₁</u>	<u>B₂</u>	<u>C₁</u>	<u>A₁</u>
a ₁	a' ₁		b ₃	b' ₃	c ₃	{a ₂ , a ₃ } X
a ₂	a' ₂					
a ₃	a' ₃					
a ₄	a' ₄					

Hence, it is not possible to have 2 relations.
and therefore we need 3 relations \Leftrightarrow minimum for
the given case

Case 3:
Many-many
+
Both sides
total
Participation.



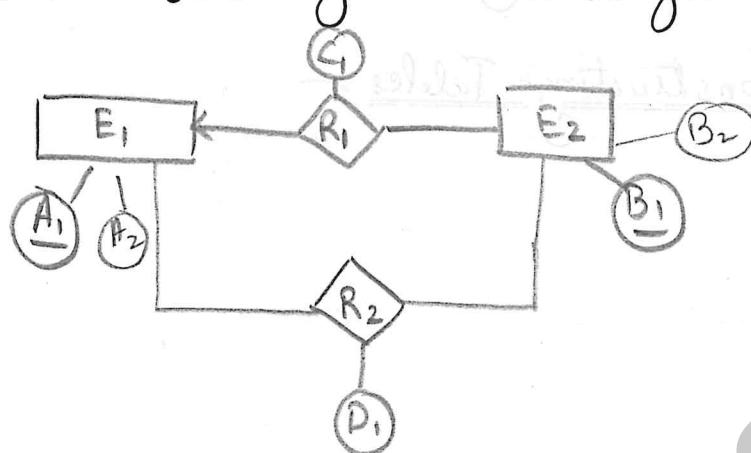
As we can see here 1 relation is possible in minimum case with $\{A_1, B_1\}$ as the primary key. But it will result in lot of redundancy of data.

<u>A₁</u>	<u>A₂</u>	<u>B₁</u>	<u>B₂</u>	<u>C₁</u>
a_1	a'_1	b_1	b'_1	c_1
a_2	a'_2	b_2	b'_2	c_2
a_3	a'_3	b_3	b'_3	c_2
a_4	a'_4	b_4	b'_4	c_1

Note:- For many-to-one, cases will remain same as one-to-many as there is only change in the perspective of taking / considering entities.

Note:- For one-to-many and many-to-one + total participation on both sides, only 1 relation is required in minimum case.

(Q1) What is the minimum number of relations required for the below given E-R diagram?

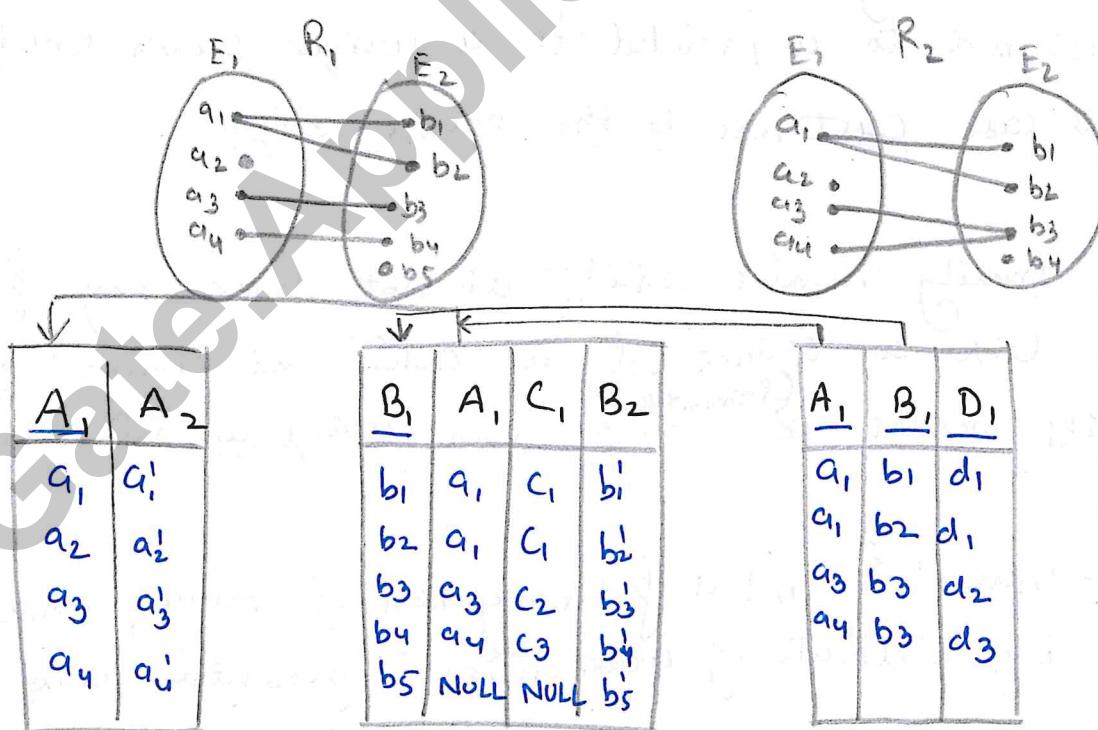


⇒ We can trivially do it in 4 relations.

~ We cannot do it in 1 relation because there is 1:m and m:n relationship with R₁ and R₂ associated respectively.

Hence, we need to check for 2 and 3 relations

→ Check for 3 relations -



Hence, minimum of 3 relations are required.

Note: we can not do it with 2 tables because there is m:n relation + Partial ∴ we need not able to get 2 tables.

III Weak and Strong Entities, self-referential relationship (2-3)

And Constructing Tables :-

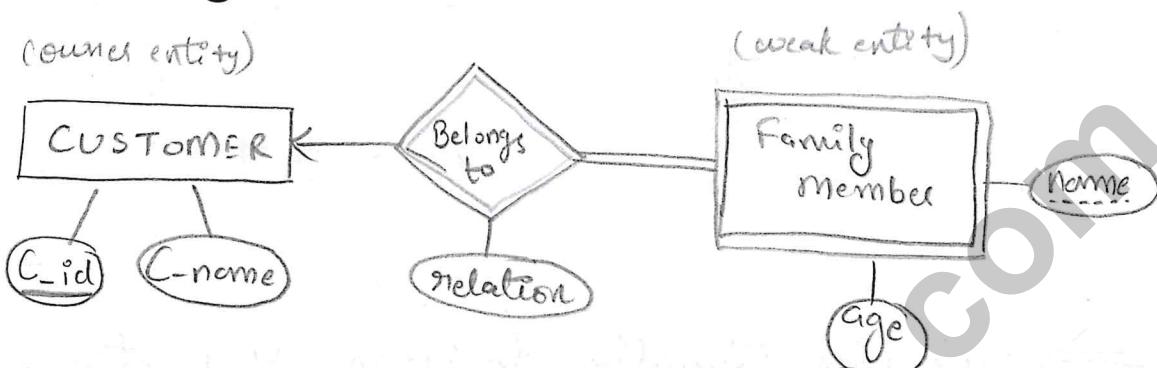


Fig: 20 Weak entity E-R diagram.

In the above E-R diagram, Family member is a weak entity. It is represented using concurrent rectangles. Note that a weak entity does not have a primary key. A weak entity cannot exist in isolation, it needs to be assigned to a / related to a unique owner entity. In this case customer is the owner entity.

Here, Family member (entity) do not have any primary key. Given a owner id ie 'C-id' and 'name' of a family member we ^(primary key) can uniquely determine a table.

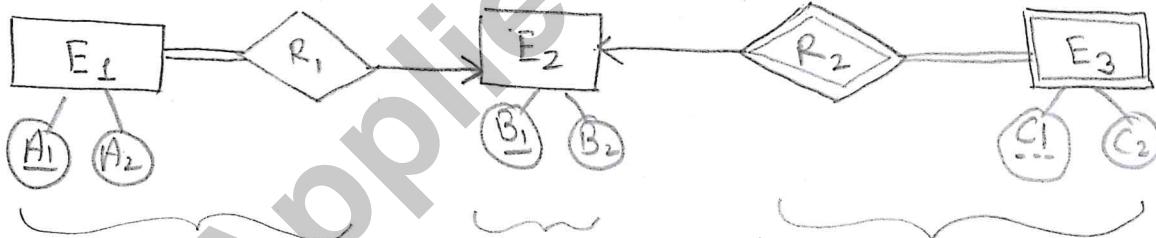
Note 'name' is not a primary key of Family member, it is a key attribute of weak entity represented using dotted line.

Note: An owner entity should always be a strong entity and relationship between owner and weak entity is always one-to-many with total participation from weak entity.

Note: In Fig: 20, the concurrent diamond / rhombus represents a weak relation between weak entity and its owner strong entity. This relationship is ^{also} called identification relationship.

Since, it is a one-to-many relationship with total participation on many side, therefore, minimum number of relation / tables required are 2. One relation for weak entity and weak relationship attributes and other is for owner entity.

(Q2) What is the minimum no. of relations req. for the given E-R diagram?



As we can see :

Fig: 21

$R_1 \rightarrow$ many-to-one with total on many side

$R_2 \rightarrow$ one-to-many with total on many side.

So, we can combine E_1 and R_1 in one relation

E_2 is the second relation

R_2 and E_3 will form the 3rd relation

So, relations will be as follows:

A ₁	A ₂	B ₁

B ₁	B ₂

B ₁	C ₁	C ₂

So, minimum 3 relations are required for above E-R diagram.

Self-referential relationships:

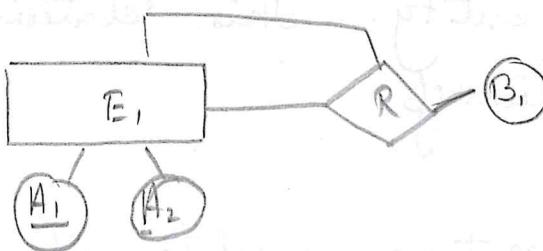
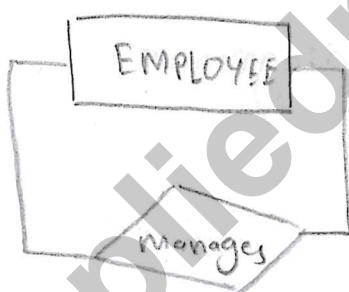


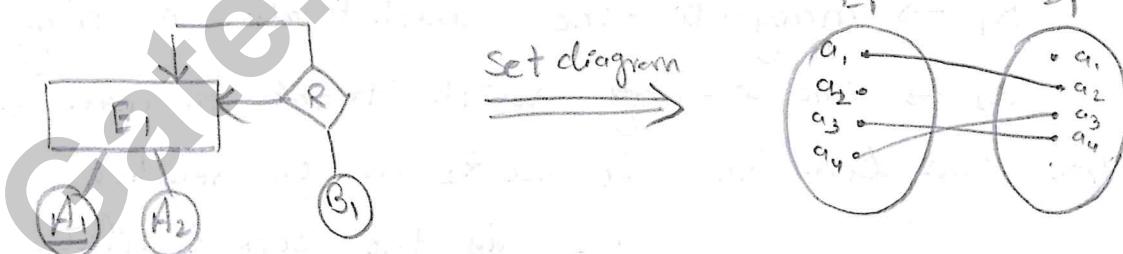
Fig: 2) Self-referential relationship

A relationship which refers to an entity and itself is called a self-referential relationship.

For example:- "An employee manages another employee" is a self-referential relationship.



Case 1: One-to-one + Partial-Partial participation.



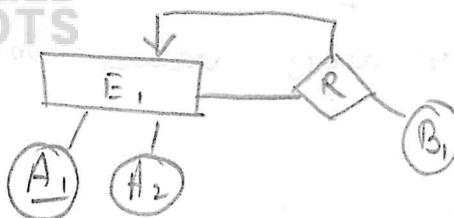
⇒ Check for 1 relation possible or not:

A ₁	A ₂	A ₁	B ₁
a ₁	a ₁ '	a ₂	b ₁
a ₂	a ₂ '	NULL	b ₁
a ₃	a ₃ '	a ₄	b ₂
a ₄	a ₄ '	a ₃	b ₃

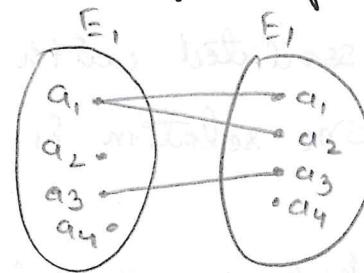
Hence, 1 relation is enough for this case. with A₁ as primary key

Phone: +91 844-844-0102

Case 2: One-to-many + Partial-partial participation:



Set diagram



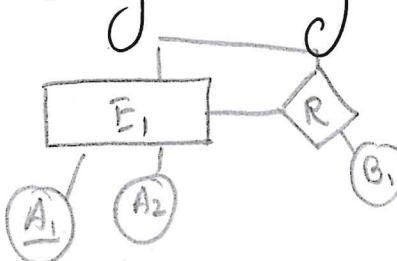
⇒ Check for 1 relation is possible or not:

refers to			
A ₁	A ₂	A' ₁	B ₁
a ₁	a' ₁	a ₁	b ₁
a ₁	a' ₂	a ₂	b ₁
a ₃	a' ₃	a ₃	b ₃
NULL	a' ₄	a ₄	NULL

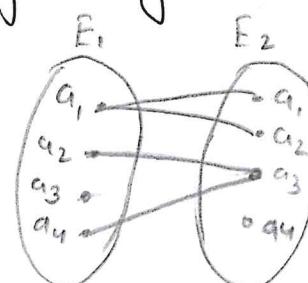
Here, we cannot make A₁ as our primary key because in set diagram we can see a₁ is related to more than one value of itself. Therefore, we will have A'₁ as our primary key.

Hence, only 1 relation is required for given condition with A'₁ as our primary key.

Case 3: Many-to-many + Partial-partial participation:



Set diagram



As we can see here ~~both~~ both A_1 and A'_1 are associated with two values or more values. Therefore, one relation is not possible.

⇒ Check for 2 relation is possible or not:

\downarrow	A_1	A_2	\downarrow	A_1	A'_1	B_1
	a_1	a_1'		a_1	a_2	b_1
	a_2	a_1'		a_1	a_2	b_1
	a_3	a_2'		a_2	a_3	b_2
	a_4	a_3'		a_4	a_3	b_3

Hence, minimum 2 relations are required.

NOTE: For self-referential relationships (any case) + total participation, always require 2 relation (minimum) (even in the case of many-to-many) but it have a lot of redundancy.

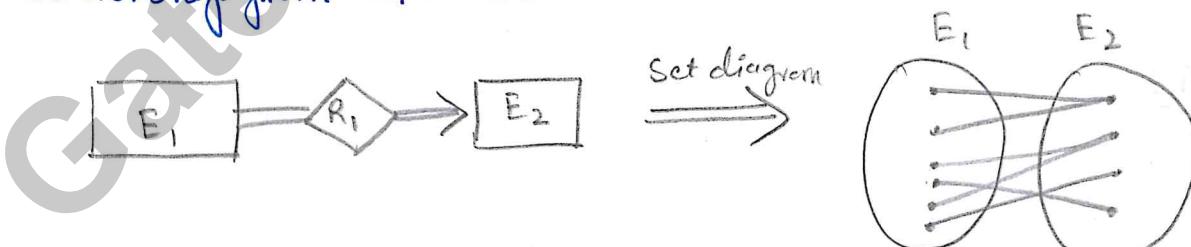
IV Solved Problems: (PYQ) Phone: +91 844-844-0102

(2.4) APPLIED ROOTS

(Q1) In an Entity-Relationship (ER) model, suppose R is a many-to-one relationship from entity set E₁ to entity set E₂. Assume that E₁ and E₂ participate totally in R and that the cardinality of E₁ is greater than the cardinality of E₂. Which one of the following is true about R?

- (A) Every entity in E₁ is associated with exactly one entity in E₂.
- (B) Some entity in E₁ is associated with more than one entity in E₂.
- (C) Every entity in E₂ is associated with exactly one entity in E₁.
- (D) Every entity in E₂ is associated with at most one entity in E₁.

⇒ It is given that R₁ is many-to-one + total participation relationship from E₁ to E₂.



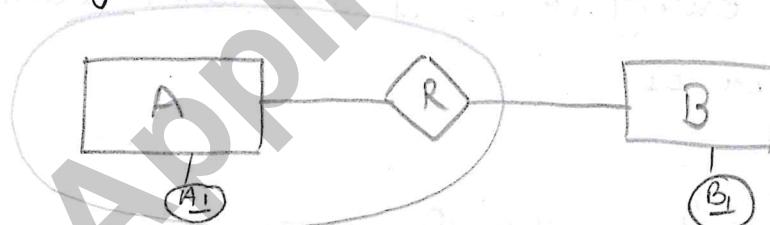
- (A) It is true, because every entity in E₁ is associated with exactly one entity in E₂ (as we can see in set diagram).
- (B) While options (B), (C), and (D) are false because according to (B) the relationship btw E₁ and E₂ is many-to-many.

According to (C) the relationship R₁ in btw E₁ and E₂ is one-to-many and option (D) is false as it is implying at most one-to-one with many-to-many btw E₁ and E₂. Hence option (A) is true.

(Q2) An E-R model of a database consists of entity types A and B. These are connected by a relationship R which does not have its own attribute. Under which of the following conditions, can the relational table for R be merged with that of A?

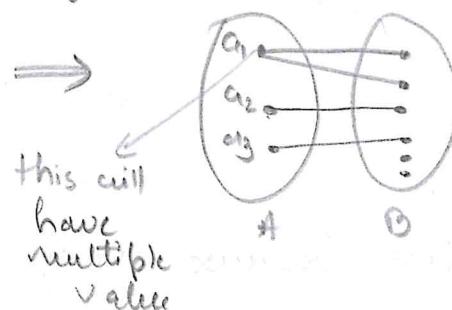
- (A) Relation R is one-to-many and the participation of A in R is total.
- (B) Relation R is one-to-many and the participation of A in R is partial.
- (C) Relation R is many-to-one and the participation of A in R is total.
- (D) Relation R is many-to-one and the participation of A in R is partial.

\Rightarrow



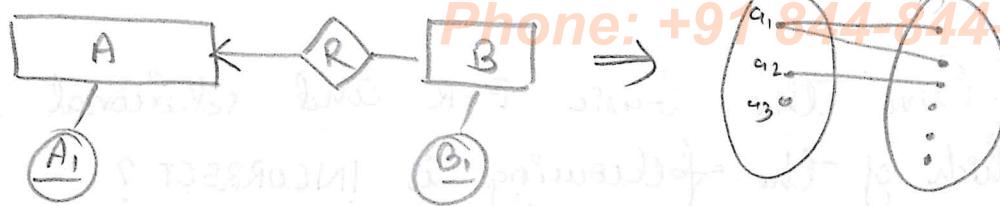
In which condition
we can merge them?

Case A:



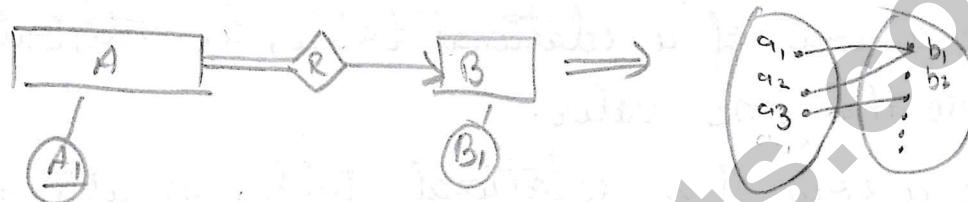
As we can see, entity A will have multiple values of 'a1'. Therefore, we cannot combine R with A in this condition. Hence, this is false.

Case B:



In case B also, there is the same problem - 'a_i' will have multiple values. Hence, R could not be merged with A.

Case C:



Here, every value in A is associated with exactly 1 value in B. Therefore R could be merged A.

Hence, this is true.

Case D:



Here, also ~~some~~ value of A is associated with exactly one value in B while some are not even related.

Hence, this is also true, But (C) is more appropriate.

Hence, answer is (C) and (D) both but (C) is more stronger option than (D) as it will not have any NULL values for R in relation. ~~as~~

(Q3) Given the basic E-R and relational models, which of the following is INCORRECT?

- (A) An attribute of an entity can have more than one value.
- (B) An attribute of an entity can be composite.
- (C) In a row of a relational table, an attribute can have more than one value.
- (D) In a row of a relational table, an attribute can have exactly one value or a NULL.

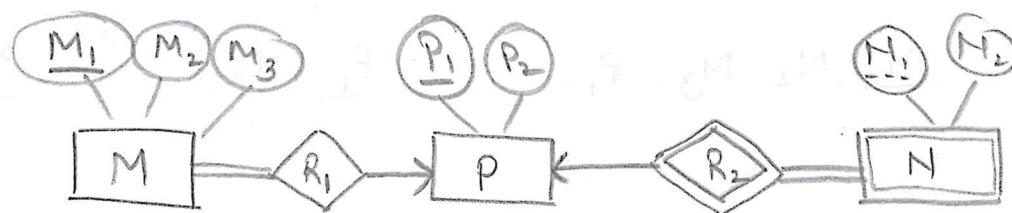


- (A) It is correct, it is talking about multivalued attributes, i.e. an attribute of an entity can have more than one value and these attributes are called multivalued attributes.
- (B) It is correct, an entity can have composite / compound attributes. Example 'name' is an attribute in 'customer' entity which is formed using 'first-name', 'last-name'.
- (C) It is incorrect, a tuple / row of a table cannot store more than ~~one~~ one value in a cell.
- (D) It is correct, it is an opposite statement of option (C).

Hence, option (C) is incorrect.

(Q4) The minimum number of tables needed to represent M, N, P, R_1, R_2 is:

- (A) 2
- (B) 3
- (C) 4
- (D) 5



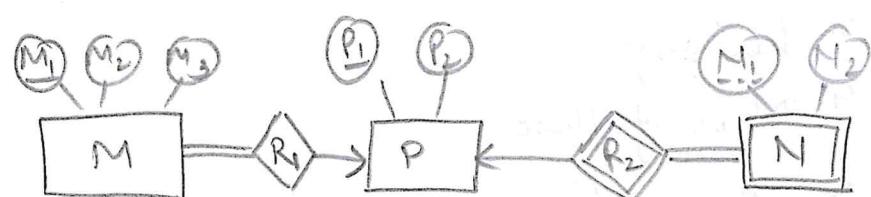
\Rightarrow If we recall, it is same/similar to Fig: 21 problem.

- option (C) and (D) are the trivial option.
- We cannot do it in 2 relations because neither M and P could be combined together nor P and N could be combined.

Therefore, option (B) is correct.

(Q5) Consider the following E-R diagram. Which of the following is a correct attribute set for one of the tables for the minimum number of tables needed to represent M, N, P, R_1, R_2 ?

- (A) M_1, M_2, M_3, P_1
- (B) M_1, P_1, N_1, N_2
- (C) M_1, P_1, N_1
- (D) M_1, P_1 .



\Rightarrow Since, it is an extension of the (Q4), we know that minimum no. of relation required are 3 : $M+R_1, P, R_2+N$

So, the 3 relations formed are:

<u>M₁</u>	M ₂	M ₃	P ₁

Table 1

<u>P₁</u>	P ₂

Table 2

<u>P₁</u>	<u>N₁</u>	N ₂

Table 3

Hence, option 1 is obviously correct as these are the attributes of table 1.

(Q6) Consider the entities "hotel-room" and "Person" with a many-to-many relationship "lodging" as shown below. If we wish to store information about rent payments to be made by person(s) occupying different hotel rooms, then this information should appear as an attribute of

1. Person
2. Hotel Room
3. Lodging
4. None of these.



⇒ Rent payment will be paid only when a "person lodged in a hotel room". (Person on his own will not pay rent and hotel-room on its own will not have rent) Only when a person stays in a hotel room then he is supposed to pay rent. Hence, there is no logic to keep attribute 'rent' in person or Hotel-Room entity. Therefore, it should be the attribute of Lodging (2)

INTRODUCTION TO RELATIONAL MODELS

I. Mathematical model of Tables:

(3.1)

Relational models or Relational databases are nothing but tables and these tables are the mathematical relations implicitly.

Relational database or models uses the same terminology. That is:-

- * Relations are also called tables.

- * Columns are also called attributes. The set of values a column can take is called domain.

- * Rows are also called tuples.

~ Overview of Relation in mathematics ~

In mathematics, a relation is defined as a subset of two sets $A \times B$

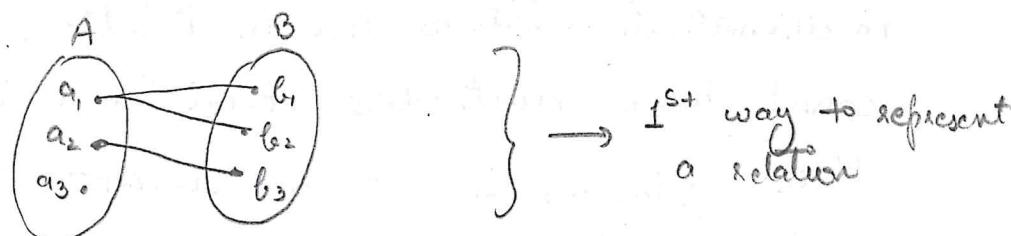
i.e. Given set $A = \{a_1, a_2, a_3, \dots, a_n\}$

$B = \{b_1, b_2, b_3, \dots, b_n\}$

$R \subseteq A \times B$

where $A \times B = \{(a_i, b_j) \mid a_i \in A, b_j \in B\}$

example: $R: A \rightarrow B$ and $R \subseteq A \times B$



So, R could be any subset of $A \times B$, ie.

$$R \subseteq A \times B \Rightarrow R = \{(a_1, b_1), (a_1, b_2), (a_2, b_3)\}$$

Mail: gatcse@appliedroots.com

2nd way to represent a relation

A	B
a ₁	b ₁
a ₁	b ₂
a ₂	b ₃

Phone: +91 844-844-0102

\rightarrow 3rd way to represent a relation (using tuple)

\rightarrow This is called a tuple in relation while in tables we call it as row.

here, domain of the attribute A is nothing but the set A which we have defined i.e. $A = \{a_1, a_2, a_3, \dots, a_n\}$. Similarly, for attribute B.

Hence, DBMS relation models and tables are related to the concept of mathematical relation and all the properties of mathematical sets and relation are applicable to DBMS relations / tables also. like

\rightarrow Sets cannot have duplicate values and same with the database models. As we have already studied that in relational model no two tuples could be exactly same. (Concept of primary key)

\rightarrow In mathematics relation $R \subseteq A \times B$ i.e. $R = \{(a_1, b_1), (a_1, \{b_1, b_2\})\}$. This relation is neither allowed in

mathematical relations nor in relations of DBMS. A tuple cannot have multivalue stored in a cell.

Here $\{b_1, b_2\}$ is not an element of B they are subset of B.

Hence, we can now prove that why there are such rules in DBMS

Relational model / E-R diagrams as at its core they are the mathematical relations / tables.

Note:- In relational database, the outline of database relationships and structure is given by relational schema. Relational schema is written as : <relation name> (<set of attributes>) + integrity constraint

e.g. $R_1 (cid, cname, czip, caddr)$ + integrity constraints

II. Solved Problems:

(3.3)

(Q1) Consider the following tables T_1 and T_2 .
 In table T_1 , P is the primary key and Q is the foreign key referencing R in table T_2 with on-delete cascade and on-update cascade. In table T_2 , R is the primary key and S is the foreign key referencing P in table T_1 with on-delete set NULL and on-update cascade. In order to delete record (3,8) from table T_1 , the number of additional records that need to be deleted from table T_2 is _____

 $T_1:$

P	Q
2	2
3	8
7	3
5	8
6	9
8	5
9	8

 $T_2:$

R	S
2	2
8	3
3	2
9	7
5	1
7	2

→ The relation T_1 and T_2 are related as follows
Phone: +91 844-844-0102

<u>on delete: Set Null</u>	
on delete: cascade	
$T_1:$	$T_2:$
P	R
2	2
3	8
7	3
5	8
6	9
8	5
9	8

This is deleted

On deleting (3,8) in T_2 , only change which will occur is in ~~the~~ table T_2 . The cell having attribute S value as 3 will become NULL.

Hence, no tuples are deleted on deleting (3,8)

Therefore, Answer is 0 (Zero)

(Note if the question how many tuples got modified then answer = 1)
(Q2) The following table has two attributes A and C where A is the primary key and C is the foreign key referencing A with on-delete cascade.

A	C
2	4
3	4
4	3
5	2
7	2
9	5
6	4

The set of all tuples that must be additionally deleted to preserve referential integrity when the tuple $(2, 4)$ is deleted is:

- (A) $(3, 4)$ and $(6, 4)$
- (B) $(5, 2)$ and $(7, 2)$
- (C) $(5, 2), (7, 2)$ and $(9, 5)$
- (D) $(3, 4), (4, 3)$ and $(6, 4)$

\Rightarrow

	A	C
①	2	4
	3	4
	4	3
②	5	2
③	7	2
④	9	5
	6	4

This is deleted initially.

This will get deleted becoz ① got deleted.

This will get deleted because of ② got deleted.

Tuple ②, ③, ④ i.e. $(5, 2), (7, 2)$ and $(9, 5)$ will get deleted on deleting $(2, 4)$.

Hence, correct option is (C).

(Q3) Let $R(a, b, c)$ and $S(d, e, f)$ be two relations in which d is the foreign key of S that refers to the primary key of R . Consider the following four operations on R and S :

I. Insert into R

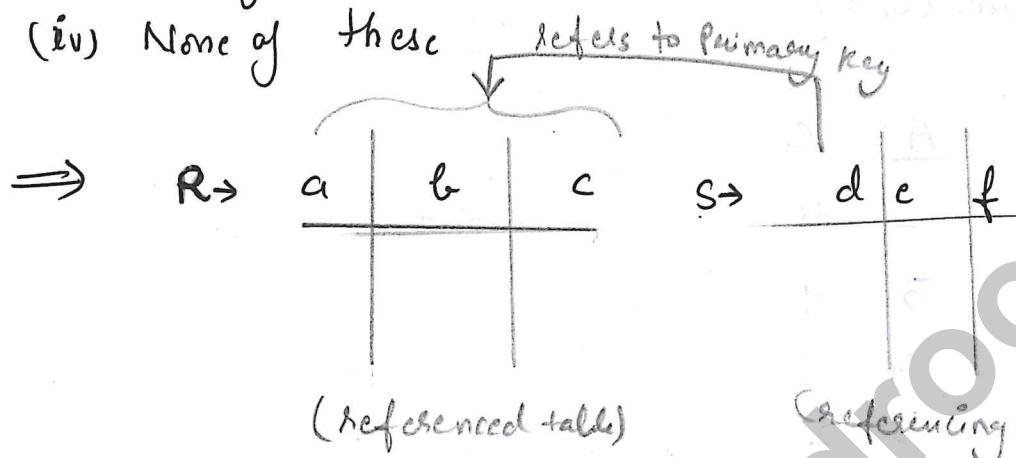
II. Insert into S

III. Delete from R

IV. Delete from S

Which of the following can cause violation of the referential integrity constraint above?

- (i) Both I and IV
- (ii) Both II and III
- (iii) All of these
- (iv) None of these



- I. Insert in R will not cause any referential integrity problem
- II. insertion in S, required to check whether inserted value of d is valid or not. Therefore, it may cause referential integrity constraint
- III. Deletion from R required tuple to either cascaded deleted or set to null. so, this may cause referential integrity problem.
- IV. Deletion from S will not cause any problem.

Hence, I and III may cause referential integrity constraint violation. So, correct option is (ii)

(Q4) The maximum number of superkeys for the relation Schema R(E, F, G, H) with E as the key is _____

⇒

$$\begin{matrix} E & U & \emptyset \\ & F & \\ & G & \\ & : & \\ & FH & \end{matrix}$$

∴ There are 3 elements left as subsets
possible = $2^3 = 8$

∴ max number of superkeys possible = 8

(Q5) Given the STUDENTS relation as shown below.

For (StudentName, StudentAge) to be the key for this instance, the value X should not be equal to _____

Student id	<u>Student Name</u>	<u>Student Email</u>	<u>Student Age</u>	CPI
2345	Shakar	Shankar@math	X	9.4
1287	Swati	Swati@ce	19	9.5
7853	Shankar	Shankar@cse	19	9.4
9876	Swati	swati@mech	18	9.3
8765	Ganesh	Ganesh@civil	19	8.7

⇒ Since, we have two entries in the name of 'Shankar'

∴ {Student Name, ~~StudentAge~~} is key. So in order to avoid these type of duplication and to uniquely determine

a row/tuple, X should not be 19

(Q6) Which of the following is NOT a superkey in a relational schema with attributes V, W, X, Y, Z and primary key VY?

1. VXYZ
2. VWXZ
3. VWXY
4. VWXYZ

⇒ A superkey is primary ~~from~~ or candidate key \cup other attributes

1. $\{ \}$ is a union of VY \cup XZ \therefore $\{ \}$ is a superkey
2. $\{ \}$ is not a superkey because it does not include the primary key.
3. $\{ \}$ is a union of VY \cup WX \therefore it is a superkey
4. $\{ \}$ is a union of VY \cup WXYZ \therefore it is a superkey

Hence, correct option is (2)

(Q7) Consider a relational table with a single record for each registered student with the following attributes:

1. Registration_Num: Unique registration number of each registered student.
2. UID: Unique identity number, unique at the national level for each citizen.
3. Bank_Account_Num: Unique account number at bank.

A student can have multiple accounts

or joint accounts. This attribute stores the primary account number.

4. Name: Name of the student
5. Hostel_Room: Room number of the hostel.

Which one of the following option is INCORRECT?

- (A) Bank_Account_Num is a candidate key
- (B) Registration_Num can be primary key
- (C) VID is candidate key if all students are from the same country
- (D) If S is a superkey such that $S \cap VID = \text{NULL}$ then $S \cup VID$ is also a superkey.



(A) Candidate key should uniquely determine a tuple.
 A student can have multiple bank account or joint accounts
 So, there is possibility that two students S_1 and S_2 are having
 joint account and they will share the Bank_Account_Num.

Hence, this is incorrect.

- (B) It is correct because Registration_Num is unique for each student.
- (C) It is correct as VID can uniquely determine a student from the same country
- (D) Suppose we have superkey say $S = \{\text{Registration_Num}, \text{Bank_Account_Num}\}$. It is able to uniquely identify a tuple. VID is one of the attribute which is not included in the Superkey as it is given in the option (ie $S \cap VID = \text{NULL}$)
 So, even if we include VID in the superkey S, than also it will be uniquely identify each tuple and is still a superkey. Hence, it is correct.

Therefore, correct option is (A)

RELATIONAL ALGEBRA

Call: +91 844-844-0102



I. Introduction to Relational Algebra and Basic Operators:

We are going to study 3 types of query languages (tools to obtain data from a DBMS):

- [1] Relational Algebra
 - [2] Relational Calculus
 - [3] SQL
- } → These are more towards mathematical view
} → This is more towards programmatic view.

→ Procedural vs Declarative languages:-

PROCEDURAL

→ Clear instructions on "how to generate outputs for given outputs" are given

→ Example: Relational Algebra, C, C++, Java. But relational algebra is not a programming language, it's a mathematical language.

DECLARATIVE

→ It tells "what we want" and don't tell anything regarding "how to get/obtain it".

→ Example: SQL

→ Basic Relational Operators: (These operators are called relational operators because they operate on relations)

→ Projection:-

It is represented as ' Π ' ($P:$)

Customer :-
Relation

Cid	Cname	Zip
i ₁	n ₁	z ₁
i ₂	n ₁	z ₂
i ₃	n ₂	z ₂

fig: 22

Projection operation is a relational operator which takes only one relation as an input (it is a unary operator) and outputs the attributes/columns listed in the form of a relation.

Syntax: $\Pi_{<\text{attribute_name}>} (<\text{Relation_name}>)$

Example 1: $\Pi_{\text{cid, cname}} (\text{Customer}) =$

↓
It has taken
a relation
as i/p

Cid	Cname
i ₁	n ₁
i ₂	n ₁
i ₃	n ₂

↓
output's
another relation
with only mentioned
fields.

Example 2: $\Pi_{\text{cname}} (\text{Customer}) =$

Cname
n ₁
n ₁
n ₂

X Wrong output.

Since, a relation is a set of tuples, therefore, it should not have duplicates. Hence, the output of example 2 query is:

Cname
n ₁
n ₂

→ Selection :-

It is denoted by symbol σ (sigma)

Selection operation is a relational operator which selects a subset of tuples satisfying the given condition.

Syntax: $\sigma_{<\text{condition}>} (<\text{relation_name}>)$

↳ It can be any propositional logic statement

Example 1: On considering the relation given in Fig: 22

the output of following query is

$\sigma_{\text{Cname} = n_1} (\text{customer}) =$

↑
Takes L
relation
as i/p

Cid	Cname	Czip
i ₁	n ₁	z ₁
i ₂	n ₁	z ₂

↑
Op's another
relation.

Example 2:

$\sigma_{(\text{Cname} = n_1) \wedge (\text{Czip} = z_1)} (\text{customer}) =$

Cid	Cname	Czip
i ₁	n ₁	z ₁

Note - We can combine ~~projection~~ and ~~selection~~ operation together.

Example:

$$\Pi_{\text{cid}} (\sigma_{\text{cname} = n_1} (\text{customer}))$$

This will get executed first

O/p of $\sigma_{\text{cname} = n_1} (\text{customer})$ is \Rightarrow

Cid	Cname	Zip
i ₁	n ₁	3 ₁
i ₂	n ₁	3 ₂

Now $\Pi_{\text{cid}} (R_1) =$ (Let's name it as R₂)

Cid
i ₁
i ₂

final o/p

Rename :-

It is represented by 'P' (rho). It is used to rename the output selection.

Syntax: P <new_name> | <old_name> (<relation_name>)

Example: Considering the relation given in fig:22, the output of the following query is

P(cust_name/cname) (customer) =

Cid	Cust_name	Zip
i ₁	n ₁	3 ₁
i ₂	n ₁	3 ₂
i ₃	n ₂	3 ₂

↓
It will only rename the attribute name.
Rest of the relation will remain

Set Operations : $(U, \cap, -)$

C_1	C_2	C_3
1	a	α
2	b	β
3	c	γ

C_1	C_2	C_3
1	a	α
4	d	δ
2	c	γ

Fig:23

All set operations are applicable to relations also because relations are the set of tuples.

NOTE- Set operation - Complement is not often used in relations. This is not often used because

$$R_1^C = \boxed{A} \quad \text{where } U = C_1 \times C_2 \times C_3$$

This condition is not often.

$$\text{In Fig:21 } R_1^C = U - R_1$$

(i) Union: (Combines all unique tuples of both R_1 and R_2 into 1 relation)
 Considering relations given in fig: 23

$$R_1 \cup R_2 =$$

C_1	C_2	C_3
1	a	α
2	b	β
3	c	γ
4	d	δ
2	c	γ

Note: (1,a, α) won't repeat because relation is a set.

(ii) Intersection: (combines all the common tuples in between 2 relations into a ~~single~~ new relation)

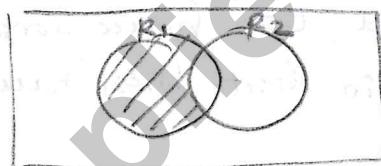
On considering relations given in fig: 23

$$R_1 \cap R_2 =$$

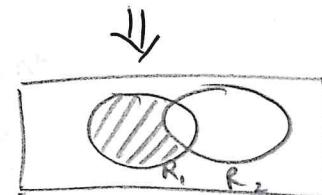
C_1	C_2	C_3
1	a	α

NOTE - Intersection is possible only when all the attributes values matches. In above example there is only one tuple whose all the three attributes are same.

(iii) Difference: (outputs all the tuples which are in relation R_1 but not in relation R_2)



$$R_1 - R_2 = R_1 - (R_1 \cap R_2) =$$



On considering the relations given in Fig: 23
the output will be:

$$R_1 - R_2 = R_1 - (R_1 \cap R_2) =$$

C_1	C_2	C_3
2	b	β
3	c	γ

NOTE: Compatibility of relations: +91 844-844-0102



We cannot perform set theoretic operations if # of attributes in R_1 and R_2 are not equal or if they are having different domain.

Ex:-

1) $R_1(C_1, C_2, C_3)$, $R_2(C_1, C_2) \rightarrow$ NOT COMPATIBLE

2) $R_1(C_1, C_2, C_3)$, $R_2(C_1, C_2, C_4) \rightarrow$ NOT COMPATIBLE
↳ Datatype is diff

here, C_3 and C_4 have different domain.

3) $R_1(C_1, C_2, C_3)$, $R_2(C_1, C_2, C_5) \rightarrow$ NOT COMPATIBLE

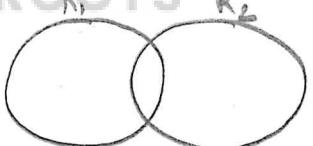
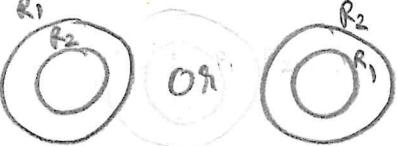
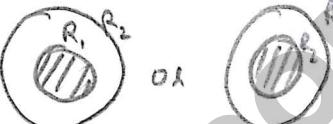
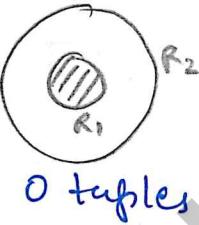
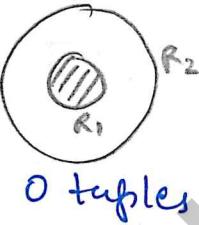
C_3 and C_5 have same domain but they have different name

4) $R_1(C_1, C_2, C_3)$, $\rho_{C_3/C_5}(R_2(C_1, C_2, C_5)) \rightarrow$ COMPATIBLE

here, C_3 and C_5 have same domain and we have renamed the C_5 attribute to C_3 .

NOTE: Given 2 relations R_1 and R_2 having m and n tuples respectively. Suppose that R_1 and R_2 are compatible and hence, we can perform set operations on them. Then, the minimum and maximum number of tuples possible on performing set operation

Mail: gatecse@appliedroots.com

SET OPERATION	MINIMUM TUPLE	MAXIMUM TUPLE
 $R_1 \cup R_2$ $(m \text{ tuple}) \quad (n \text{ tuple})$	 $\max(m, n)$ tuples	 $m+n$ tuples
 $R_1 \cap R_2$	 0 tuples	 $\min(m, n)$ tuples
 $R_1 - R_2$ $(R_1 - R_2 \neq R_2 - R_1)$	 0 tuples	 m tuples

→ ← Cross-product / Cartesian product:

C ₁	C ₂	C ₃
1	a	x
2	b	y

C ₃	C ₄
α	β
γ	α

Fig: 24

Cartesian product is helpful to merge columns from two relations. It is given as $R_1 \times R_2$.

$R_1 \times R_2$

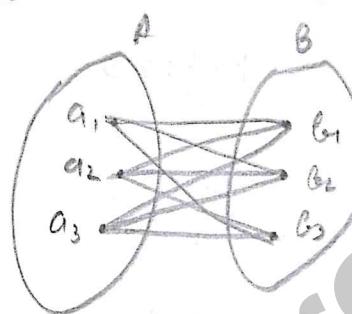
=

C_1	C_2	C_3	C_4	C_5
1	a	x	α	β
1	a	x	1	α
2	b	y	α	β
2	b	y	1	α

Phone: +91 844-844-0102

2x2 = tuples

If R_1 has m tuples and R_2 has n tuples, then $R_1 \times R_2$ will have $m \times n$ tuples. It can be understood as:

NOTE:

→ Suppose a relation R_1 with m attributes and R_2 with n attributes

$$R_1(r_1, r_2, \dots, r_m) \times R_2(s_1, s_2, \dots, s_n)$$

↓ attribute

$$= \{ (r_1, r_2, \dots, r_m, s_1, s_2, \dots, s_n) \mid \begin{array}{l} (r_1, r_2, \dots, r_m) \in R_1 \\ \text{AND} \\ (s_1, s_2, \dots, s_n) \in R_2 \end{array} \}$$

$m+n$ attributes.

→ We can do cross product of any relation with any number of relations.

Eg. $R_1 \times R_1$, $R_1 \times R_2 \times R_3$, $R_1 \times R_1 \times R_1$, etc all are valid cartesian Product

→ Note that, duplicate rows are not allowed but duplicate columns are allowed.

→ ↵ Derived Operators

Phone: +91 844-844-0102

**APPLIED
ROOTS**

These operations are derived using basic operators.
Ex Natural join, division etc.

C_1	C_2	C_3
1	a	α
2	b	β

R

C_2	C_3	C_4
a	α	1
a	γ	2

S

Fig: 25

(i) Natural join:

It is represented using ' \bowtie '. It is a binary operator which takes 2 relations as I/P and outputs a new relation. It is similar to cross product but with a twist. The twist is, there should be atleast one column which is common in between the two relations in order to join them. (all other columns should be compatible i.e. having same domain)

On considering relation given in fig: 25

$$R \bowtie S = \left[\begin{array}{|c|c|c|c|} \hline C_1 & C_2 & C_3 & C_4 \\ \hline 1 & a & \alpha & 1 \\ \hline \end{array} \right]$$

Here only one set of (C_2, C_3) in R got matched with set of (C_2, C_3) in S. Hence, we have only one tuple in resultant relation.

It could be written as:

$$R \bowtie S = \overline{\cap}_{\substack{R.C_2 = S.C_2 \\ R.C_3 = S.C_3}} (R \times S) \quad (\text{it is expanded expression of natural join using basic operators})$$

This will select a bunch of columns.

This will do filtering of tuples obtained in cartesian product

This will give us 6 columns i.e. C₁, C₂, C₃, C₄, C₅, C₆. So, we will use R.C₂ and R.C₃ in order to avoid confusion.

(ii) Conditional join:

It is a slight modification of natural join. It is also referred as "theta join". It is same as natural join with a condition associated with it. It is written as ' \bowtie_{θ} ' (θ -join) or ' \bowtie_c ' (c stands for condition)

$$R \bowtie_c S = \overline{\cap}_c (R \times S)$$

For given relation R and S in fig: 25,

$$R \bowtie_{R.C_1 < S.C_4} S = \overline{\cap}_{R.C_1 < S.C_4} (R \times S)$$

with these conditions

R.C ₁	R.C ₂	R.C ₃	S.C ₂	S.C ₃	S.C ₄
1	a	α	a	γ	2

Note: Conditional join do not have projection operation

Mail: gates@appliedroots.com This is one of the difference between Conditional join and theta join. Another difference btw them is

The conditional join can have any condition which could be defined using operators $<$, $>$, \leq , \geq etc. While natural join can have only one condition which is ' $=$ ' and that too between common attributes or fields of two relations R and S.

These are the fundamental difference between natural join and conditional join.

(ii) Left outer join:

It is denoted as ' \bowtie '. This operation allows us to keep all tuples in the left relation. If there is no matching tuple found in right relation, then the attributes of right relation in the join result are filled with NULL values.

It could be written as

$$R \bowtie S = R \bowtie S \cup \text{all tuples from the left relation that failed the join with NULL values appended}$$

On Considering relation R and S in fig 25

$$R \bowtie S =$$

C_1	C_2	C_3	C_4
1	a	α	1
2	b	β	NULL

Fig: 26

→ matched tuple
→ unmatched tuple with appended NULL

NOTE - Left outer join could be both conditional and non-conditional. If it is conditional, then condition will be applied at selection part of natural join.

It may look like:

$$R \bowtie_c S = R \bowtie_c S \cup \{ \text{all tuples from the left relation that failed join with NULL values appended} \}$$

NOTE- Using basic operators, it could be written as:

$$R \bowtie S = R \bowtie S \cup (R - \Pi_{r_1, r_2, \dots, r_m}(R \bowtie S)) \times (\underbrace{\text{NULL, NULL}, \dots, \text{NULL, NULL}}_{\text{No. of NULLs}})$$

where, $R = (r_1, r_2, \dots, r_m)$
 $S = (s_1, s_2, \dots, s_n)$

No. of NULLs =
 No. of additional column for R.

Here, in our example there is only one additional column i.e 'C₄' so we will do cross product of only one NULL in the given case.

(iv) Right outer join:-

It is denoted as ' \bowtie ' . This operation allows us to keep all tuples in the right relation. If there is no matching tuple found in left relation, then the attributes of left relation in the join result are filled with NULL values.

It could be written as :

$$R \bowtie S = R \bowtie S \cup \text{all tuples from the right relation that failed the join with NULL values appended}$$

Everything else is same as that of left outer join.

Mail: gatcse@appliedroots.com

Using basic operators it could be written as **Phone: +91 844-844-0102**

$$R \bowtie S = R \bowtie S \cup (NULL, NULL \dots) \times (S -$$



$$\Pi_{s_1, s_2 \dots s_n} (R \bowtie S)$$

NULL is equal to no. of additional column for S

On considering relation R and S given in fig: 25

$$R \bowtie S =$$

C ₁	C ₂	C ₃	C ₄
1	a	α	1
NULL	a	γ	2

Fig: 27

(V) Full outer join:

It is represented as ' \bowtie '. It could be given as

$$R \bowtie S = R \bowtie S \cup R \bowtie S$$

C ₁	C ₂	C ₃	C ₄
1	a	α	1
2	b	β	NULL
NULL	a	γ	2

} it is union of all tuples obtained in $R \bowtie S$ and $R \bowtie S$

(Vi) Division:

Cust-id	Category
1	C ₁
1	C ₂
2	C ₁
3	C ₂
4	C ₁
4	C ₂

R

Category
C ₁
C ₂

S

Fig: 28

Suppose we are asked a query/question that - (from fig:28)

"Customers who purchased from every category"

→ According to the relations given in fig: 28, there are only two categories available ~~and~~ (C₁ and C₂) and there are only 2 customers who purchased product from both of them. Hence, the output will be

Cid
1
4

T

Note that we are able to arrive to this solution with the help of division operation ($R \div S$) also written as ' R/S '.

$R \div S$ or $R/S \Rightarrow$ It outputs those tuples from R, for which it holds that all their combinations with tuples in S are present in R. Here, only 1 and 4 have their

Combination with all tuples in S (and are present in R) - 0102

$$\therefore R \div S = T$$

now, if we do $T \times S =$

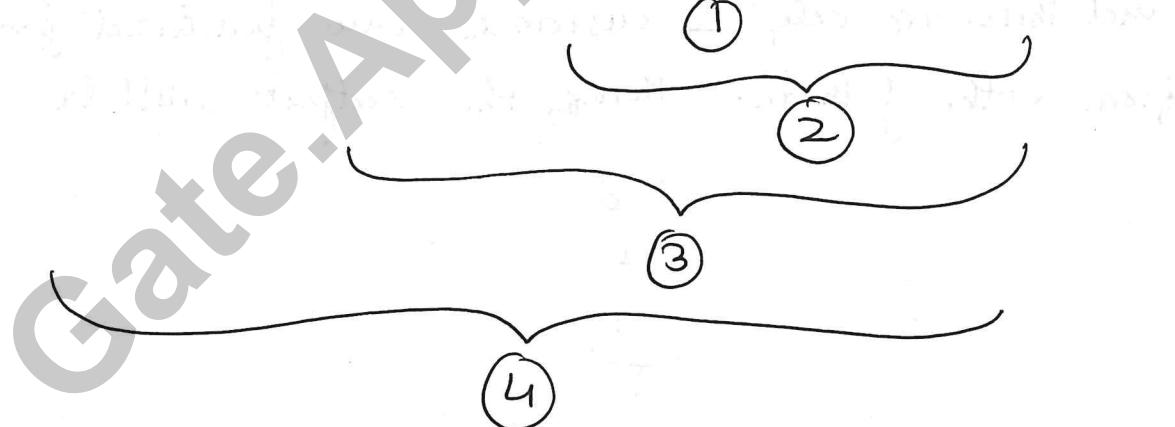
Cust-id	category
1	C ₁
1	C ₂
2	C ₁
2	C ₂

→ As we can see here, $R \supseteq T \times S$

We can write division operator using basic ~~attribute~~ operators as:

$$R(\text{Cust-id}, \text{Category}) \div S(\text{Category}) =$$

$$\overline{T}_{\text{Cust-id}}(R) - \overline{T}_{\text{Cust-id}}\left(\overbrace{\left(\overline{\Pi}_{\text{Cust-id}}(R) \times S\right)}^1 - R\right)$$



Output of ① :

Cust-id	Category
1	C ₁
1	C ₂
2	C ₁
2	C ₂
3	C ₁
3	C ₂
4	C ₁
4	C ₂

Output of (2) :-

cust-id	category
2	C ₁
3	C ₂

Phone: +91 844-844-0102

Output of (3) :-

cust-id
2
3

Output of (4) :-

cust-id
1
2
3
4

custid
2
3

custid
1
4

NOTE - Completeness :- We can represent any relational operator using { U, -, X, Π , σ } in relational algebra. This is a very popular complete set of operators in relational algebra.

II Solved Problems :-

(4.3)

(Q1) Consider the following relation P(x, y, z), Q(x, y, t) and R(y, v);

How many tuples will be returned by the following relational algebra query?

$$\Pi_x (\sigma_{(P.y = Q.y \wedge R.v = v_2)} (P \times R)) - \Pi_x (\sigma_{(Q.y = R.y \wedge Q.t > 2)} (Q \times R))$$

X	Y	Z
x_1	y_1	z_1
P		
x_1	y_1	z_2
x_2	y_2	z_2
x_2	y_4	z_4

X	Y	T
x_2	y_1	2
Q		
x_1	y_2	5
x_1	y_1	6
x_3	y_3	1

Y	V
y_1	v_1
R	
y_3	v_2
y_2	v_3
y_2	v_2

⇒ The query given is:

$$\Pi_x (\sigma_{(P.Y = Q.Y \wedge R.V = V_2)} (P \times R)) - \Pi_x (\sigma_{(Q.Y = R.Y \wedge Q.T > 2)} (Q \times R))$$

①

②

- ① will return only one tuple because of the 2 conditions $P.Y = Q.Y$ and $R.V = V_2$ hence, output will be:

X
x_2

- ② will return also return only one tuple because of two conditions $Q.Y = R.Y$ and $Q.T > 2$

hence, the output will be:

X
x_1

final we need to find $\textcircled{1} - \textcircled{2} = \textcircled{1} - (\textcircled{1} \cap \textcircled{2}) \Rightarrow$

Mail: gatecse@appliedroots.com $\textcircled{1} \cap \textcircled{2} = \text{NULL}$

$\therefore \textcircled{1} - \textcircled{2} = \text{m}$ - Hence only 1 tuple will be returned.

(Q2) Consider the relation $r(A, B)$ and $s(B, C)$, where $S \cdot B$ is a primary key and $r \cdot B$ is a foreign key referencing $S \cdot B$.

Consider the query:

$$Q : r \bowtie (\sigma_{B < 5}(s))$$

Let LOJ denote ^{natural} left outer join operation. Assume that r and s contain no NULL values. Which one of the following is NOT equivalent to Q ?

- (1) $\sigma_{B < 5}(r \bowtie s)$
- (2) $\sigma_{B < 5}(r \text{ LOJ } s)$
- (3) $r \text{ LOJ } (\sigma_{B < 5}(s))$
- (4) $\sigma_{B < 5}(r) \text{ LOJ } s$

\Rightarrow We are given that:

A	B
α	4
β	5

B	C
4	a
5	b

r s

Query given is: $r \bowtie (\sigma_{B < 5}(s)) \Rightarrow Q \Rightarrow$

A	B	C
α	4	a

Let us consider consider an example tree(mentioned above)

(1) the output of 1st query will be:

A	B	C
α	4	a

(2) the output of 2nd query will be:

A	B	C
α	4	a

(3) the output of 3'd query will be: [+91 9448440102](tel:+919448440102)



A	B	C
a	b	c
B	S	NULL

Hence, (3) is the answer ($\epsilon_3 \neq \emptyset$)

(Q3) Consider a database that has the relation schema CR (StudentName, CourseName). An instance of the schema CR is as given below:

The following query is made on the database:

$$T_1 \leftarrow \Pi_{\text{CourseName}} (\sigma_{\text{StudentName} = 'SA'} (CR))$$

$$T_2 \leftarrow CR \div T_1$$

The number of rows in T_2 is _____

StudentName	CourseName
SA	CA
SA	CB
SA	CC
SB	CB
SB	CC
SC	CA
SC	CB
SC	CC
SD	CA
SD	CB
SD	CC
SD	CD

SE	CD
SE	CA
SE	CB
SF	CA
SF	CB
SF	CC

CR

→ Here T_1 is $\Rightarrow \text{π}_{\text{CourseName}}(\sigma_{\text{StudentName} = 'SA'}(R))$

This will
o/p all
tuples with
name = 'SA'
it will project
course name of all
student having name
as 'SA'

CourseName
CA
CB
CC

T_2 is $\Rightarrow CR \div T_1$ } look at all StudentName which
are combined with {CA, CB, CC}

StudentName
SA
SC
SD
SF

} all these 4 are
combined with
{CA, CB, CC}

Mail: gatcse@appliedroots.com

Hence No. of tuples = 4 //

(Q4) Consider two relations $R_1(A, B)$ with tuples $(1, 5), (3, 7)$ and $R_2(A, C) = (1, 7), (4, 9)$. Assume that $R(A, B, C)$ is the full natural outer join of R_1 and R_2 . Consider the following tuples of the form (A, B, C) : $a = (1, 5, \text{NULL})$, $b = (1, \text{NULL}, 7)$, $c = (3, \text{NULL}, 9)$, $d = (4, 7, \text{NULL})$, $e = (1, 5, 7)$, $f = (3, 7, \text{NULL})$, $g = (4, \text{NULL}, 9)$. Which one of the following statements is correct?

- (1) R contains a, b, c, f, g but not c, d
- (2) R contains all a, b, c, d, e, f, g
- (3) R contains e, f, g but not a, b
- (4) R contains e but not f, g .

\Rightarrow here $R_1 \Rightarrow$

A	B
1	5
3	7

A	C
1	7
4	9

$$R = R_1 \bowtie R_2$$

A	B	C
1	5	7
3	7	NULL
4	NULL	9

} actual tuples in R .

As we can see in the option only option (3) and (4) are having 3 and 1 tuple respectively. Therefore, these 2 are

potential options.

As we can observe

$$c = (1, 5, 7)$$

$$f = (3, 7, \text{NULL})$$

$$g = (4, \text{NULL}, 9)$$

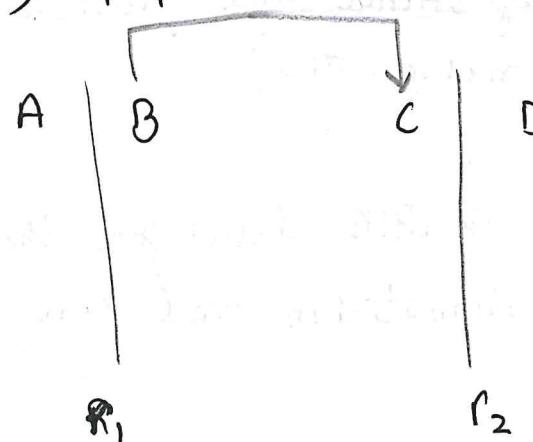
Since, option (3) includes all of these 3 tuples
hence correct option is (3).

(Q 5) Suppose (A, B) and (C, D) are two relation schemas. Let r_1 and r_2 be the corresponding relational instances. B is a foreign key that refers to C in r_2 . If data in r_1 and r_2 satisfy referential integrity constraints. Which of the following is ALWAYS TRUE?

- (A) $\Pi_B(r_1) - \Pi_C(r_2) = \phi$
- (B) $\Pi_C(r_2) - \Pi_B(r_1) = \phi$
- (C) $\Pi_B(r_1) = \Pi_C(r_2)$
- (D) $\Pi_B(r_1) - \Pi_C(r_2) \neq \phi$

\Rightarrow

Given:



all values in
 B should be
present in C

According to the options ~~Phone: +91 844-0102~~

option (A) is appropriate answer as it will be always true. ($B - C = \emptyset$)

5. TUPLE RELATIONAL CALCULUS

I (S.1) Tuple Relational calculus - I :-

- * TRC is a declarative language. Here we only mention "what we want and not how to get it"
- * TRC is a mathematical foundation for SQL.
- * TRC is heavily based on propositional and predicate logic.
(Note: Propositional and predicate logic is covered in discrete mathematics)
- * Edgar Codd @ IBM developed the TRC, RA, SQL. He is a mathematician and computer scientist.

A query in TRC looks like **Phone: +91 844-844-0102**

$$\{ t \mid P(t) \}$$

↑
it refers
to a
tuple

↓
it is a predicate
logic expression.
involving operation like
 $\{ \vee, \wedge, \neg, \exists, \forall, \in \}$

Cid	Cname	Czip
i ₁	n ₁	z ₁
i ₂	n ₁	z ₂
i ₃	n ₂	z ₃
i ₄	n ₃	z ₁

Customers

Fig: 29

Let us consider this example relation to understand TRC with the help of RA

→ ↵ Select rows from customer where zip code = z₁

In RA $\Rightarrow \sigma_{Czip=z_1} (Customer)$ (it is an example of Selection operation)

↔ Both will give same op

In TRC $\Rightarrow \{ t \mid t \in customer \wedge t.Czip = z_1 \}$

↪ Predicate logic expression

It will be read as: "Return me the set of tuples t, such that t belongs to customer table and t.Czip = z₁"

NOTE:

There are many ways of writing TRC query. For example, in place of writing "t.Czip" we can also write "t[Czip]".

Both are equivalent. In place of writing "t ∈ customer"

We can also write it as "Customer(t)" So, there are multiple notations of writing TRC.

→↓ Project cid from customer relation

In RA $\Rightarrow \Pi_{cid}(\text{customer})$

In TRC $\Rightarrow \{t \mid \exists s (\text{se customer} \wedge s.cid = t.cid)\}$ more popular notation

It will be read as: "Set of all tuples such that there exists S, such that S belongs to customer table and S.cid = t.cid."

This seems to be very round about way to write as we are creating a new variable S.

Another way we can write this in TRC is

$\{t.cid \mid t \in \text{customer}\}$ \rightarrow This is one of the way to write but it is not a popular way.

NOTE -

Let's break $\{t \mid \exists s (\text{se customer} \wedge s.cid = t.cid)\}$

There exist S

S belongs to customer means S is a tuple now.

"it is a tuple" it has all 3 columns

"we want only cid that is why we use t.cid only" only cid attribute is used here. It will return only cid attribute as op

NOTE- In previous example of projection, variable t is not associated with any \exists and \forall while s is associated with \exists and \forall , therefore, variable s is called Bounded variable. t is unbounded or free variable. t is not bounded by any first order logic or predicate logic condition (\exists and \forall)

→ Project cid and $cname$ from customer

In RA $\Rightarrow \Pi_{cid, cname}(\text{customer})$

In TRC $\Rightarrow \{t \mid \exists s \in \text{customer} (s.cid = t.cid \wedge s.cname = t.cname)\}$

It involves $P(t.cid, t.cname)$

∴ it will print only $t.cid$ and $t.cname$ and not $t.zip$

→ Project cid for all $czip=3$, on customer

(Selection + Projection)

In RA $\Rightarrow \Pi_{cid} (\sigma_{czip=3}(\text{customer}))$

In TRC $\Rightarrow \{t \mid \exists s \in \text{customer} (s.czip = 3 \wedge s.cid = t.cid)\}$

cid

Selection operation

Projection operation

It will be read as: "Select set of all tuples t such that there exist s such that belongs to customer and $s.czip = 3$, and $s.cid = t.cid$ "

NOTE Same like RA, In TRC **Also we can write complete phone: +91 948410102**

Set of operations ie $\{U, -, \pi, \sigma, X\}$ and hence, using these basic operators we can write all derived operators like join, divisions etc in TRC.

→ Cartesian product

Given relations:

$\{\text{Customer}(cid, cname, zip)\}$
 $\{\text{Product}(cid, pid, time)\}$

In RA \Rightarrow Customer X Product

(it will have 6 attributes)

In TRC $\Rightarrow \{t \mid \exists t_1 \in \text{Customer}, \exists t_2 \in \text{Product} \quad \begin{array}{l} t_1 \text{ is row of customer} \\ t_2 \text{ is row of product} \\ \text{every row of } t_1 \text{ is} \\ \text{combined with every } t_2 \text{ row.} \end{array}\}$

$t \cdot cid = t_1 \cdot cid \wedge$
 $t \cdot cname = t_1 \cdot cname \wedge$
 $t \cdot zip = t_1 \cdot zip \wedge$
 $t \cdot cid_2 = t_2 \cdot cid \wedge$
 $t \cdot pid = t_2 \cdot pid \wedge$
 $t \cdot time = t_2 \cdot time\}$

It will be read as: "Set of all tuples t such that
t₁ belongs to C and t₂ belongs to P (C for customer
and P for product) such that {set of all conditions}"

→ Set operations.

(i) Union

If two relations R and S are set operation compatible
than

In RA $\Rightarrow R \cup S$

In TRC $\Rightarrow \{t \mid t \in R \vee t \in S\}$

OR

$\{t \mid R(t) \vee S(t)\}$

(ii) Set difference:



In RA $\rightarrow R - S$

t should not be equal to
any tuple in S (ie $t \notin S$)

In TRC $\Rightarrow \{t \mid t \in R \text{ AND } t \notin S \text{ (} t_1 \neq t_2 \text{)} \}$

It can also be written as $\neg(t = t_1)$

It will be read as "for all tuples set t , where t belongs to R AND for all t_1 belongs to S , t is not equal to t_1 "

* find customer id of customers who made atleast one purchase.

In RA $\Rightarrow Tl_{cid} (\sigma_{c.cid = p.cid} (customer \times product))$

Selection

Cross Product

In TRC $\Rightarrow \{t \mid \exists t_1 \in \text{Customer}, \exists t_2 \in \text{Product} (t_1.cid = t_2.cid \wedge t.cid = t_1.cid) \}$

Projection

\hookrightarrow will return $t.cid$

$P(t.cid)$

It will be read as "A set of tuples t , there exists t_1 belongs to customer and t_2 belongs to product such that $t_1.cid$ should be equal to $t_2.cid$ and $t.cid = t_1.cid$ "

NOTE-

Formula / equations in predicate logic

$$1) \neg(\forall x P(x)) = \exists x (\neg P(x)) \quad \left. \begin{array}{l} \text{Demorgan's laws in} \\ \text{Predicate logic} \end{array} \right\}$$

Similarly, $\neg(\exists x P(x)) = \forall x (\neg P(x))$

$$2) P \rightarrow Q \equiv \neg P \vee Q$$

$$3) \forall x (P(x) \wedge Q(x)) = \forall x P(x) \wedge \forall x Q(x)$$

4) $\forall u (P(x) \rightarrow Q(u)) = \forall u (\neg P(u) \vee Q(u))$

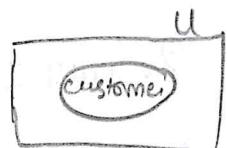
II (S.2) Safe queries & Domain Relational calculus:

→ Safe vs unsafe Queries:

Let's take an example to understand this:

$$\{t \mid \underbrace{\neg(t \in \text{customer})}_{\substack{\text{This means } t \notin \text{customer} \\ \text{OK}}}\} \rightarrow \exists t \text{ means, I want all tuples that do not belong to customer.}$$

This is a valid statement in predicate logic. But if we talk about this in relational calculus, it is nothing but compliment of relation customer.



In set theory compliment is given as :

$$\text{customers}^c = U - \text{customers} \quad (U \text{ stands for compliment})$$

$$= \text{infinite result.}$$

So, the TRC statement written above is a valid statement but the resultant set is a non-finite set of tuples. ($\exists t$ is like every human being possible which is not a customer of Amazon). Therefore, it is an unsafe query.

Hence, any TRC query which results in a non-finite set of tuples is called a unsafe query.

NOTE - Unsafe TRC queries could not be written in Relational algebra. (In relational algebra it is already assured that there is no concept of complement). For every safe TRC query there is an equivalent Relational algebra query.

The expressive power of Relational algebra is same as that of safe TRC. Unsafe TRC is a valid TRC statement query but the result is non-finite.

NOTE -

There is a trick to determine whether given TRC is safe or unsafe. If there is a complement of free variable then most likely it is going to be unsafe query.

$$\{t \mid \neg(t \in \text{customer})\}$$

↓
free variable

→ Domain Relational calculus (DRC):

RA, TRC and DRC are mathematical foundation or model. These are not used while writing a query in real world. It is used to optimise the query and not to work with them.

→ Selection: (For relation given in fig:29)

In RA \Rightarrow ~~t1 t2~~ $\sigma_{c_2:p=3}(\text{customer})$

In TRC \Rightarrow $\{t \mid t \in \text{Customer} \wedge t.c_2:p=3\}$

↓
tuple variable

$\exists n \text{ DRC} \Rightarrow \{ \langle u_1, u_2, u_3 \rangle \mid \underbrace{\langle u_1, u_2, u_3 \rangle}_{\text{Domain variables}} \in \text{customer} \}$

u_1 can take all values that belongs to the domain cid.

u_2 can take all values that belongs to the domain cname. Similarly for u_3 .

In TRC we take whole tuple / row as variable while here every attribute is an variable

(u_1) cid	(u_2) cname	(u_3) cgb

customer
 $\langle u_1, u_2, u_3 \rangle$ is an ordered triplet

NOTE -

Syntax of DRC:

$\{ \langle u_1, u_2, \dots, u_n \rangle \mid \underbrace{P(u_1, u_2, \dots, u_n)}_{\text{Predicate logic on these domain variables.}} \}$

So only difference between DRC and TRC is that in TRC we have ~~tuple~~ variable for a Tuple while in DRC we have variables for all attributes.

\Rightarrow Projection:

$\exists n \text{ RA} \Rightarrow \Pi_{cid} (\text{customer})$

$\exists n \text{ DRC} \Rightarrow \{ \underbrace{\langle u_1 \rangle}_{\substack{\downarrow \\ \text{this will select only } \\ \text{cid as } \\ u_1 \text{ is for } \\ \text{cid}}} \mid \exists u_2, u_3 (\langle u_1, u_2, u_3 \rangle \in \text{customer}) \}$

$\exists n \text{ TRC} \Rightarrow \{ t \mid \exists s (s \in \text{customer} \wedge s.cid = t.cid) \}$

Mail: gatcse@appliedroots.com

→ Cartesian product: Given **Phone: +91 844-844-0102**
 Customer (cid, channel, zip)
 Product (pid, bid, time)

In RA \Rightarrow Cartesian product Customer \times Product

In DRC $\Rightarrow \{ \langle u_1, u_2, u_3, u_4, u_5, u_6 \rangle \mid \exists (u_1, u_2, u_3) \in \text{Customer}$
 $\wedge \exists (u_4, u_5, u_6) \in \text{Product} \}$

III (S.3) Solved problems:

(Q1) Which of the following relational query languages have the same expressive power?

- I. Relational Algebra
- II. TRC restricted to safe expressions
- III DRC restricted to safe expressions

- (A) II and III only
- (B) I and II only
- (C) I and III only
- (D) I, II and III

→ As we know RA is equivalent to TRC and DRC

(Condition that TRC and DRC are restricted to safe expression)

Hence, all the 3 statements are equivalent.

Therefore, correct option is (D)

(Q2) Which of the following relational calculus expressions is not safe?

- $\{ t \mid \exists u \in R_1 (t[A] = u[A]) \wedge \neg \exists s \in R_2 (t[A] = s[A]) \}$
- $\{ t \mid \forall u \in R_1 (u[A] = "x") \Rightarrow \exists s \in R_2 (t[A] = s[A] \wedge s[A] = u[A]) \}$
- $\{ t \mid \neg (t \in R_1) \}$
- $\{ t \mid \exists u \in R_1 (t[A] = u[A]) \wedge \exists s \in R_2 (t[A] = s[A]) \}$

\Rightarrow While evaluating all the query we can see (a), (b) and (d) are very complex and we can spend our time to evaluate them but since we can see option (c) is definitely a unsafe query (as it is giving the complement as output). Therefore, we can ~~even~~ say the (c) is the correct option.

(Q3) Consider the relation Employee (name, sex, SupervisorName) with name as the key. SupervisorName gives the name of the supervisor of the employee under consideration. What does the following tuple Relational Calculus query produce?

$$\{ e.name \mid \text{employee}(e) \wedge (\forall x) [\neg \text{employee}(x) \vee x.\text{SupervisorName} \neq e.name \vee x.sex = "male"] \}$$

- (A) Names of employees with a male supervisor
 (B) Names of employees with no immediate male subordinate
 (C) Names of employees with no immediate female subordinate
 (D) Names of employees with a female supervisor.

\Rightarrow Given query

{ e.name | employee(e) }
 it will return names of employee

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

We will try to convert this negation by "for all" expression to "there exists" because we are not comfortable with the negations.

$$\therefore \forall u \{P(u)\} = \exists u \neg P(u) = \exists u \neg \neg P(u)$$

So it could be written as:

{ e.name | employee(e) }

$\exists (\exists x) (\text{employee}(u) \wedge u.\text{supervisorName} = e.\text{name} \wedge$
 $u.\text{sex} = \text{"female"}]$

This whole expression is now negated.

So, it will be read as: "I need employee name, ~~such that~~ there

does not exist x , such that x is an employee and x is supervisor and is female". In short we want e.name such that

(Q4) Which of the following is/are equivalent to $\forall t \in r(P(t))$? Applied Roots 841 14-0102

- I. $\neg \exists t \notin r(P(t))$
- II. $\exists t \notin r(\neg P(t))$
- III. $\neg \exists t \in r(\neg P(t))$
- IV. $\exists t \in r(\neg P(t))$

- (A) I only
 (B) II only
 (C) III only
 (D) III and IV only

⇒ Here:

$r \rightarrow$ relation

$P(t) \rightarrow$ predicate logic

$\forall t$ could also be written as

$$\forall t \in r(P(t)) = \underbrace{\forall t (t \in r(P(t)))}_{\text{More appropriate}}$$

$$= \neg (\neg \forall t (t \in r(P(t))))$$

$$= \neg \exists t \underbrace{(t \notin r(P(t)))}_{\text{This is not in optim}}$$

↓

This could be written as

$$t \notin r(P(t)) = t \in r(\neg P(t))$$

Mail: gatcse@appliedroots.com

hence, correct option is (c),

$$\neg \exists t \in r(\neg P(t)) \Rightarrow \text{III}$$

STANDARD QUERY LANGUAGE (SQL)

- * SQL stands for standard Query language
- * SQL is a standard way to query / obtain / add / delete / modify data
- * It is used by data scientist, machine learning engineer, S/w engineer etc
- * SQL is not a general purpose programming language. It is a domain specific language whose primary task is to build database. Other languages like C/C++/Java/Python are general purpose programming language and therefore use to build web servers, operating system, while SQL is not used for such tasks.
- * SQL is called declarative programming language where we do not specify step-by-step procedure to get / obtain data.

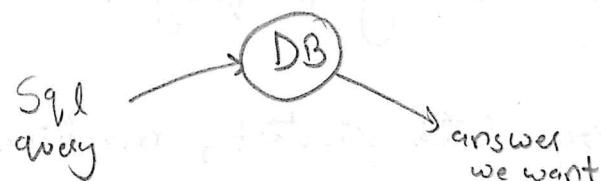
Ex- SELECT movie_name from movies WHERE Year > 200
 here, we are declaring "what we want?" and
 not, "how to get it?"

Because it is declarative, that is why it becomes simple ~~for the~~ to use.

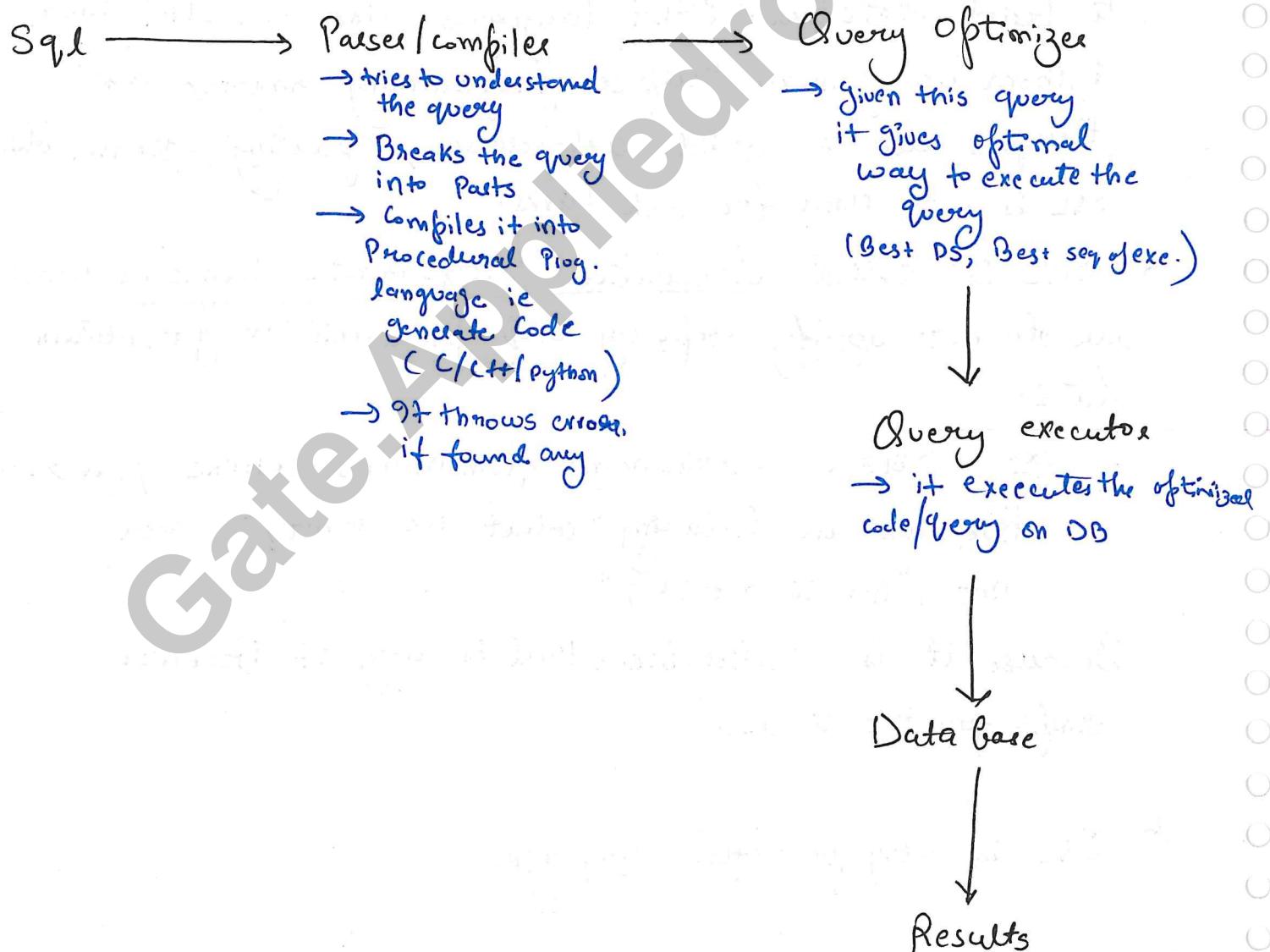
- * SQL is very powerful language.

7.3 Execution of SQL statement

Given the fact that SQL is an declarative programming language, "how does an SQL statement is executed?"

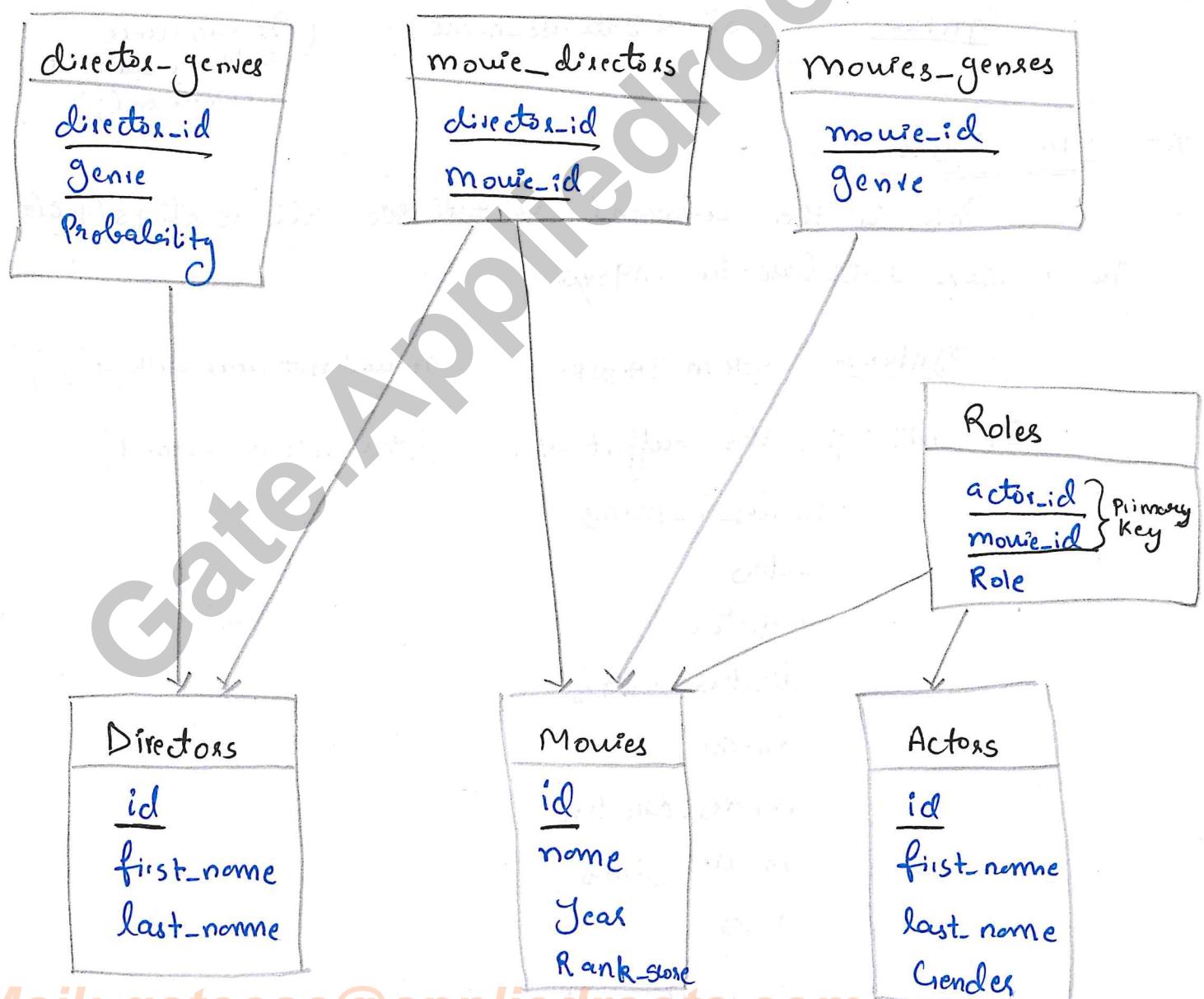


So, internally it works as follows:



7.4 IMDB Dataset:

- * IMDB is a website owned by Amazon. This dataset will be considered as an set example to understand different SQL commands. This dataset is a movie dataset.
 - * This data set have 388,269 movies from 1888-2008 which have 817,718 actors and 86,880 directors.
- Structure of the IMDB dataset (Relational database):



SQl COMMANDS

P.1 Use, Describe , Show-tables:

USE :

It is used to tell mysql that which database we want use or work on.

Syntax - USE <Database_name>; (It can work with or without semicolon (;))

SHOW TABLES :

This is the command to see all relations/tables in a given database in mysql .

Syntax - SHOW TABLES; (It will not work without ';')

It will give the output as :

(for IMDB dataset)

Tables_in_IMDB
actors
directors
directors_genres
movies
movies_directors
movies_genres
roles

DESCRIBE:

To understand each of this table, we use a command DESCRIBE. It will give complete info about a relation like , data type, default values, keys etc.

Syntax- DESCRIBE <Table-name>; (will not work without Semicolon (;))

Example- DESCRIBE actors;

output \Rightarrow

Field	Type	Null	Key	default	Comments
id	int(11)	No	PRI	0	
first-name	Varchar(100)	YES	MUL	NULL	
last-name	Varchar(100)	YES	MUL	NULL	
Gender	Char(1)	YES		NULL	

datatype

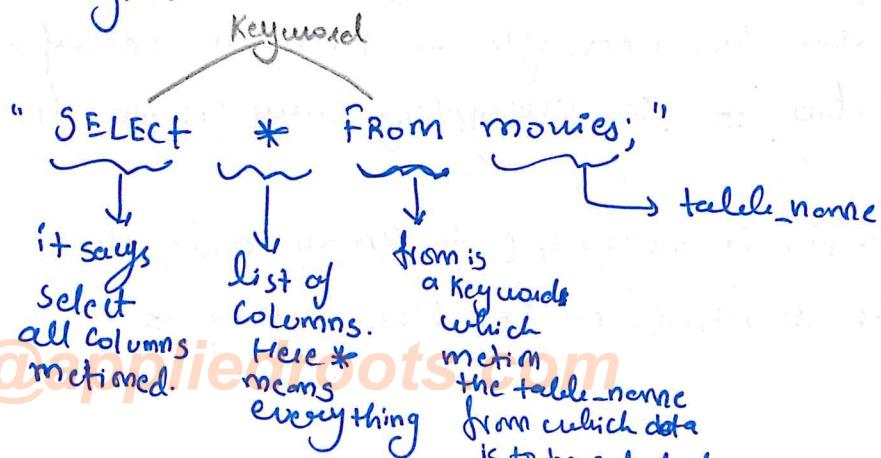
4 rows in set (0.00 sec)

→ it means
more than
1 actor could
have same first-name
or last-name.
Ex:- Both Salman
and Shahrukh
have last name
'Khan'

Q.2 Select:

Suppose we are given with a task : "list all movies in the database".

Given a query "SELECT * FROM movies;"



It will be read as "Select ~~Everything~~ + (*) from the table name movies". So, as soon as we run it, it will give complete table content as the result.

NOTE - ~~Select~~ movies will only result in the description of every attribute present in table movies. While, select basically ~~says~~ says which column to select and will give all values / content present under that column.

Example 2: "SELECT name, year FROM movies;"

→ it will give 388,205 row as o/p ~~with~~ with only two attributes "name and year" from table movies.

NOTE - Here we are not explaining how to generate output, we only mention what we want.

NOTE - The output generated by "Select" in SQL query is called "result set". It is nothing but another table that it creates. (It is a set of rows with column name)

NOTE - Example 2 is faster than ~~example 1~~ example 1 (list all movies) bcoz amount of data to be transferred is less in example 2 than in example 1 (∴ it is selecting *, ie everything / every column from table movies).

So, it is suggested to query only those attributes which are required rather than using *.

Phone: +91 844-844-0102

NOTE: The order of column name could be anything. It need not to be in the same order as it is given in the table/relation. The o/p or result set will have the order of the attributes as given in the query.

Note that, the order of rows in result set will be same as order in the original table/relation. This is called row-order preservation (So, row orders are preserved by simple select queries, there are some queries where row order is not preserved)

8.3 Limit, Offset:

Limit is a keyword which tells the number of rows we need to limit in our result set.

Example: SELECT name, year FROM movies LIMIT 20;

	Name	Year
1	Movie 1	Year 1
2	Movie 2	Year 2
3	Movie 3	Year 3
4	Movie 4	Year 4
5	Movie 5	Year 5
6	Movie 6	Year 6
7	Movie 7	Year 7
8	Movie 8	Year 8
9	Movie 9	Year 9
10	Movie 10	Year 10
11	Movie 11	Year 11
12	Movie 12	Year 12
13	Movie 13	Year 13
14	Movie 14	Year 14
15	Movie 15	Year 15
16	Movie 16	Year 16
17	Movie 17	Year 17
18	Movie 18	Year 18
19	Movie 19	Year 19
20	Movie 20	Year 20

↳ # of stories

it will give
only first
20 rows as
outputs.

Here, we got the first 20 rows only. What if we want to get the next 20 rows now. For that we use another keyword called OFFSET.

When we add offset <value> it means "ignore the ^{starting} rows equal to offset value and print the next rows equal to limit value."

Example 2: SELECT name, year FROM movies +91 844-844-0120

O/P → it will ignore first 20 rows and will print next 20 rows of table movies with columns name and year.

So, offset basically says number of rows from start to ignore. If offset is equal to 40, it will ignore first 40 rows of table movies.

Limit and offset makes the data more readable as they provide us subset of data.

Q.4 Order by:

This keyword is used to get the ordered output for a given column name.

If the order type is not mentioned, it by default order the data in ascending order (ASC)

Example 1: SELECT name, year FROM movies ORDER BY year DESC LIMIT 10;

The result
is ordered
based on
Year

The result is
in descending
order.

The feature like "sort by" and "filter" on various websites run the SQL command with ORDER BY at

Mail: gatcse@appliedroots.com

Phone: +91 811-944-102

Hence, it helps us to sort the data and therefore the result set will not be same as the original table rows order.

Q.5 Distinct:

This keyword is used to get unique / distinct values from a column / attribute.

Example- SELECT DISTINCT genre FROM movies_genres;

↑
Keyword ↓
Column-name

It will o/p all distinct genres. The o/p set will have 21 rows for 21 different genres for the given IMDB dataset.

If we want to find distinct over multiple fields we can write it as:

Example 2: SELECT DISTINCT first_name, last_name
FROM Directors; ORDER BY first_name;

Distinct will work/applied on both of these.

It will give 86844 rows as o/p, there are two rows with first_name = "Zvonko" because there are two "Zvonko" with different last name i.e. 'Coh' and 'Saksida' respectively.

So, when distinct is applied over multiple columns then complete row should be distinct (one column value can repeat but should have different value in other columns)

8.6 Where, Comparison operators, NULL: **Phone: +91 844-844-0102**



WHERE:

This keyword is used to filter the result and apply conditions.

Example 1: `SELECT name, year, rank-score FROM movies WHERE rank-score > 9;`

 condition

 it says generate
only those rows
which have rank-score > 9

and finally for the rows which have `rank-score > 9` we will select these `name, year` and `rank-score`.

It will output 1069 rows which are stated > 9 .

Example 2: `SELECT name, year, rank-score FROM movies WHERE rank-score > 9 ORDER BY rank-score DESC LIMIT 20;`

\rightarrow It will give top 20 movies, where `rank-score > 9`.

NOTE: Whatever comes after `WHERE` clause / command is called 'Condition'.

So, a condition's output could be True, False and NULL and it may have comparison operators like $=$, $\underline{<} \neq$, $<$, \leq , $>$, \geq

not equal to

Example 3:

SELECT * FROM movies_genres WHERE
genre = "Comedy";

→ It will select all columns where genre is equal to comedy. It will o/p 56425 rows.

Example 4:

SELECT * FROM movies_genres WHERE
genre <> "Horror";

OR

SELECT * FROM movies_genres WHERE
genre != "Horror";

→ It will select all columns ~~where~~ of all genres except "horror".

→ NULL:

NULL in SQL means that : Value is unknown / missing / does not exist

NULL is a special keyword in SQL. For any attribute value as NULL means that value does not exists.

NOTE

The operator " = " does not work with NULL

Example 5: SELECT name, year, rankscore FROM movies WHERE rankscore = NULL;

Mail: gatees@appliedroots.com
⇒ O/P will be an empty set. We can't compare NULL with "equal to" symbol.

In order to make it work, we write the query using "IS" and "IS NOT".

Example 2 - SELECT name, year, rank-score FROM movies WHERE rank-score IS NULL; LIMIT 20;



it will work to give all rows where rank-score is NULL.

- It will op 20 rows with rank-score as NULL

Example 3 - SELECT name, year, rank-score FROM movies WHERE rank-score IS NOT NULL LIMIT 20;

will op all rows where rank-score is not having NULL value.

8.7 logical operators:

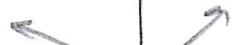
AND, OR, NOT, ALL, ANY, BETWEEN, EXISTS, IN, LIKE, SOME are some of the logical operators in SQL. Logical operators are required whenever we need to put multiple filters.

AND

Example 1 : SELECT name, year, rank-score FROM movies WHERE rank-score > 9 AND year > 2000;

Condition 1

Condition 2



These 2 Conditions are combined using AND operator.

→ It will list all movies that are released after year 2000 and having rank-score > 9. It will o/p 250 rows.

NOT

Example 2: SELECT name, year, rank-score FROM movies
WHERE NOT year <= 2000 LIMIT 20;

this Condition
is equivalent to
 $year > 2000$

→ It will list all movies that are released after year 2000 and limit the result set to the top 20 rows.

OR

Example 3: SELECT name, year, rank-score FROM movies
WHERE rank-score > 9 OR year > 2007;

Condition 1

Condition 2

↓
Combined Condition 1 and 2
and will result to TRUE if
any of the condition gets true.

→ It will list all movies that are either released after 2007 or having rank-score > 9.

NOTE - ALL and ANY are discussed in topic Subqueries.

BETWEEN

Example 4: SELECT name, year, rank score FROM movies
 (with between) WHERE year BETWEEN 1999 AND 2000;

(without between) SELECT name, year, rank score FROM movies.
 WHERE year \geq 1999 AND year \leq 2000;

\Rightarrow It will list all movies released in year between 1999 to 2000.

This is called inclusive range because both the values i.e. 1999 and 2000 are included in result.

Example 5: SELECT name, year, rank score FROM movies
 WHERE year BETWEEN 2000 AND 1999

\Rightarrow This query will not work because low value is greater than high value

NOTE

low value should be \leq high value otherwise we will get an empty set as result.

Phone: +91 844-844-0102

IN

Example 6: SELECT director_id, genre FROM directors_genres
 (with IN) WHERE genre IN ('Comedy', 'Horror');

OR

(without IN) SELECT director_id, genre FROM directors_genres
 WHERE genre = 'Comedy' OR genre = 'Horror';

⇒ It will list all director_id and genre where genre is either Comedy or horror. So, IN clause allows us to test if expression matches any value present in the list or set of values.

LIKE

(Used for text matching)

Example 7: SELECT name, year, rank_score FROM movies
 WHERE name LIKE 'Tis%';

(here % is the wildcard character to imply zero or more ~~than~~ one character)

⇒ It will list all movies name, year, rank_score where names start with 'Tis' (ie 'Tis' is followed by anything or zero)

(In shell command 'grep' or regular expression in ToC works like this)

Example 8: SELECT first-name, last-name FROM actors WHERE first-name LIKE '%es';

⇒ This will o/p all first-name, last-name of actors whose first-name are ending with 'es'

Example 9: SELECT first-name, last-name, FROM actors WHERE first-name LIKE '%es%';

⇒ This will o/p all first-name, last-name of actors whose first-name contains substring 'es'

Example 10: SELECT first-name, last-name FROM actors WHERE first-name LIKE 'Agn_s';

('_' (underline) is also a wildcard character which implies exactly one character)

⇒ This will o/p all first-name, last-name of actors whose first-name is like 'Agn_s'

↑
only 1
character
can fix
in to this.

NOTE

Consider the below relation:

T :

Name	percentage
n ₁	59%
n ₂	63%
n ₃	96%

Query: SELECT * FROM T
WHERE Percentage = '96%'

As we know % is a wildcard character which means zero or more ~~one~~ symbol.

While here % is not a wildcard character, it is a symbol. In such situations, in order to use '%' or '_' (underscore) as symbol and not as wildcard character we use backslash (\) before the '%' or '_'. Phone: +91 844-844-0102

Example - SELECT * FROM T WHERE percentage = "96\%" ↓
This will be interpreted as symbol and not as wildcard.

Note: backslash (\) : is called /interpreted as escape character in SQL.

2.8 Aggregate functions: COUNT, MIN, MAX, AVG, SUM =

These functions compute a single value on a set of rows and returns the aggregate value.

MIN

Example :- SELECT MIN(year) FROM movies;
#1937

aggregate function Column name on which aggregate function is applied

⇒ It will return a single value which is the minimum year value in table movies.

MAX

Phone: +91 844-844-0102

APPLIED
ROOTS

Example 2: SELECT MAX(Year) FROM movies;

⇒ It will return a single value which is the maximum value of year in column year in table movies.

SUM

Example 3: SELECT SUM(Year) FROM movies;

⇒ It will return a single value which is sum of all values in column year.

Similarly, avg function will find the average of all values given in a column.

So, sum and avg are the numerical operators while min and max are comparison operators.

COUNT(*)

Example 4: SELECT COUNT(*) FROM movies;

⇒ Count(*) will consider all columns and count the number of rows. It will not look for distinct rows, it will just count each and every row. Hence, output will be 388269 rows.

NOTE - Count(*) consider the rows containing NULL values also.

Mail: gatecse@appliedroots.com

Example 5: SELECT COUNT (year) FROM movies;

→ It will consider column year and will count every row without considering ~~#~~ of distinct rows.

It will output the same no. of rows as we got in Example 4 ie 388269 rows.

NOTE - Except COUNT(*), none of the aggregate function considers NULL values.

8.9 Group by :

This command is used to collect data across multiple records and group the results by one or more column.

Example 6: SELECT year, COUNT (year) FROM movies

GROUP BY year;

⇒ It will create a table with attribute year and Count all the rows having release of that year. i.e all the rows having the same year are grouped together. (grouping of all rows having some year). Here COUNT (year) will work on the group of year formed by GROUP BY clause. COUNT (year) actually gives no. of rows in that year group.

Example 7: SELECT year, COUNT (year) FROM movies

GROUP BY year ORDER BY year;

→ Everything will be same except result will be sorted in asc. order

Example 3: SELECT year, COUNT(year) year_count FROM movies GROUP BY year ORDER BY year_count;

→ It will generate some o/p but the rows will be ordered in ascending order of count(year) or year_count (year_count is an alias). It will show the year as first row which has minimum releases (as it is ordered in ascending order) and last row of o/p will have maximum releases.

NOTE- Here we have used alias because we cannot use aggregate functions in Order by clause.

NOTE- Group_by are often used with COUNT, MIN, MAX or sum.

NOTE- If grouping column contain NULL values, all null values are grouped together.

8.10 Having

This command is often used in combination with the Group_by clause to restrict the groups of returned rows to only those whose the condition is TRUE.

Example 1: SELECT year, COUNT(year) year_count FROM movies GROUP BY year HAVING year_count > 1000;

→ Here, first we are grouping and from that we will get year_count and now we will print only those year groups whose year_count > 1000 i.e. only those years which have > 1000

Here, order of execution will be :-

Step 1: GROUP BY year to create groups

Step 2: Now apply COUNT(year) or any aggregate function

Step 3: finally, apply the HAVING condition.

NOTE - Typically HAVING is used in combination with GROUP BY, but it is not mandatory.

Example - SELECT name, year FROM movies
HAVING year > 2000;



SELECT name, year FROM movies
WHERE year > 2000;

Hence, HAVING clause without GROUP BY is same as OR works same as WHERE clause.

Example - SELECT year, COUNT(year) year_count FROM movies WHERE rank_score > 9 GROUP BY year HAVING year_count > 20;

① applied over individual rows of movies table

⇒ here, in this query we have both WHERE and having clause. The seq order of execution is numbered above where is applied on individual rows while having is applied on top of groups. If we don't have group by, having will also be applied on individual rows.

NOTE - Having is applied after grouping while where is used before grouping.

→ This query will return **Phone +91 8448440102**
Such that, that year has more than 20 movies and
each of them have rank-score > 9.

Q.11 Order of Keywords:

While writing a sql query the order of
Keywords will be as follows:

- [1] SELECT
- [2] FROM
- [3] WHERE
- [4] GROUP BY
- [5] HAVING
- [6] ORDER BY (contains ASC and DESC after it)
- [7] LIMIT

9.1 Join and Natural Join:

Join's are used to combine data in multiple tables/relations.

Example: SELECT m.name, g.genre FROM movies m
 JOIN movies_genres g ON m.id = g.movie_id
 LIMIT 20;

alias of movie ↑
 alias of movie_genres ↓
 condition based on which m and g are joined

On considering fig: 31, we can see that relation movies and movies_genres are given as:-

Movies m:

Movies	
id	Name
1	n ₁
2	n ₂
3	n ₃

Genre g:

movies_genres	
movie_id	genre
1	g ₁
1	g ₂

refer to					
id	name	year	rank-score	movie_id	genre
1	n ₁	:	:	1	g ₁
2	n ₂	:	:	1	g ₂
3	n ₃	:	:	2	g ₃
				3	g ₂
				1	:

⇒ It will point name and genre

by combining table movies m and genre g
 where m.id = g.movie_id

Natural join:

A join where we have the same column-names across two tables.

example

$T_1:$

C_1	C_2

$T_2:$

C_1	C_2	C_3	C_4

Select * FROM T_1 JOIN T_2 ;

→ It will automatically joint T_1 and T_2 based on C_1 . It is equivalent to saying?

SELECT * FROM T_1 JOIN T_2 ON $T_1.C_1 = T_2.C_1$;

We can also write it as:

SELECT * FROM T_1 JOIN T_2 USING (C_1);

It will return C_1, C_2, C_3, C_4 in the result set.

9.2 Inner, left, right and Outer joins:

Consider the below relation set:

$T_1:$

C_1	C_2	C_3
1	a	b
2	c	d
4	e	f
5	g	h

$T_2:$

C_4	C_5	C_6
1	0	1
2	1	0
3	0	1
5	1	0

$\rightarrow \leftarrow$ Inner join:

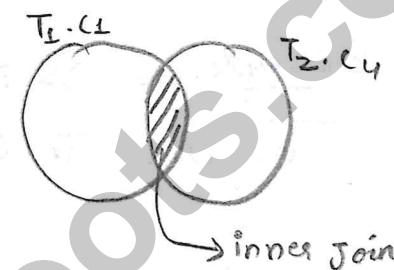
Example 1: SELECT * FROM T_1 JOIN T_2 ON $T_1.C_1 = T_2.C_4$

\Rightarrow op will be

$T \Rightarrow$

C_1	C_2	C_3	C_4	C_5	C_6
1	a	b	1	0	1
2	c	d	2	1	0
3	g	h	5	1	0

Inner join diagrammatically :-



$\rightarrow \leftarrow$ Left + outer join:

Example 2: SELECT * FROM T_1 LEFT OUTER JOIN T_2
ON $T_1.C_1 = T_2.C_4$

\Rightarrow op will be

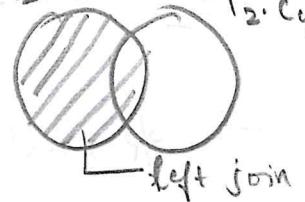
$T \Rightarrow$

No match
for this
in right relation
hence, we have
added NULL for C_4, C_5, C_6

C_1	C_2	C_3	C_4	C_5	C_6
1	a	b	1	0	1
2	c	d	2	1	0
3	e	f	NULL	NULL	NULL
4	g	h	5	1	0

here, we are combining the left relation with right relation if there are matching rows and if no matching rows are found the left relation row is clubbed with NULL

Left outer join diagrammatically \Rightarrow



NOTE - We can either use keyword LEFT OUTER JOIN or LEFT JOIN, both are same.

\Rightarrow Right Outer join:

Example 3: SELECT * FROM T₁ RIGHT OUTER JOIN T₂ ON T₁.C₁ = T₂.C₄

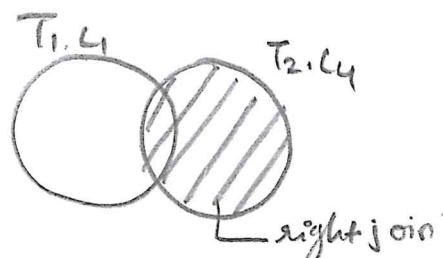
\Rightarrow o/p will be:

T ₁	C ₁	C ₂	C ₃	C ₄	C ₅	C ₆
1	a	b	f	g	0	1
2	c	d	2	l	0	0
5	g	h	5	1	0	0
Null	Null	Null	3	0	1	0

In This row eight relation do not have

T₁.C₁ = T₂.C₄ hence, we have filled C₁, C₂, C₃ with NULL

Right outer join diagrammatically \Rightarrow



NOTE - We can either use keyword RIGHT OUTER JOIN or RIGHT JOIN, both are same.

→ Full outer join:

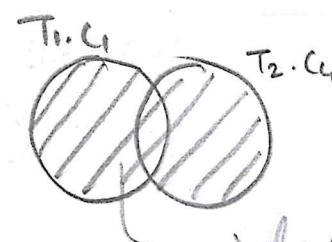
Example 4: SELECT * from T₁ FULL OUTER JOIN T₂ ON T₁.C₁ = T₂.C₄

⇒ o/p will be:

Left join ∪
Right join

C ₁	C ₂	C ₃	C ₄	C ₅	C ₆
1	a	b	1	0	1
2	c	d	2	1	0
4	e	f	NULL	NULL	NULL
5	g	h	5	1	0
NULL	NULL	NULL	3	0	1

Full outer join diagrammatically ⇒



NOTE- We can use either FULL OUTER JOIN or FULL JOIN while writing a query, both are same.

NOTE- INNER JOIN ≡ JOIN (Both are same)
(inner is optional keyword here)

NOTE- We can have 3-way join (Joining 3 relations) and K-way joins(Joining k relations)

Example- SELECT a.first_name, a.last_name FROM actors a
JOIN roles r ON a.id=r.actor_id JOIN movies m ON
m.id=r.movie_id WHERE a.id=11

APPLIED ROOTS This query is combining 3 relations ~~in~~ in a SQL query

NOTE - Joins are the most expensive computationally when we have large tables. It is one of the famous area of research that how to reduce its computation time.

SUB QUERIES / NESTED QUERIES / INNER QUERIES

Syntax:

SELECT <column-names> FROM <table-names>
WHERE <column-name> OPERATOR
(SELECT <column-name> FROM
<table-name> WHERE);

OR

SELECT Column-name [, Column-name]
FROM Table1 [, Table2]
WHERE Column-name OPERATOR (SELECT
Column-name [, Column-name]
FROM Table1 [, Table2] [WHERE]);

IN

Example: list all actors in the movies Schindler's list

⇒ Note that schindler's list is the whole series

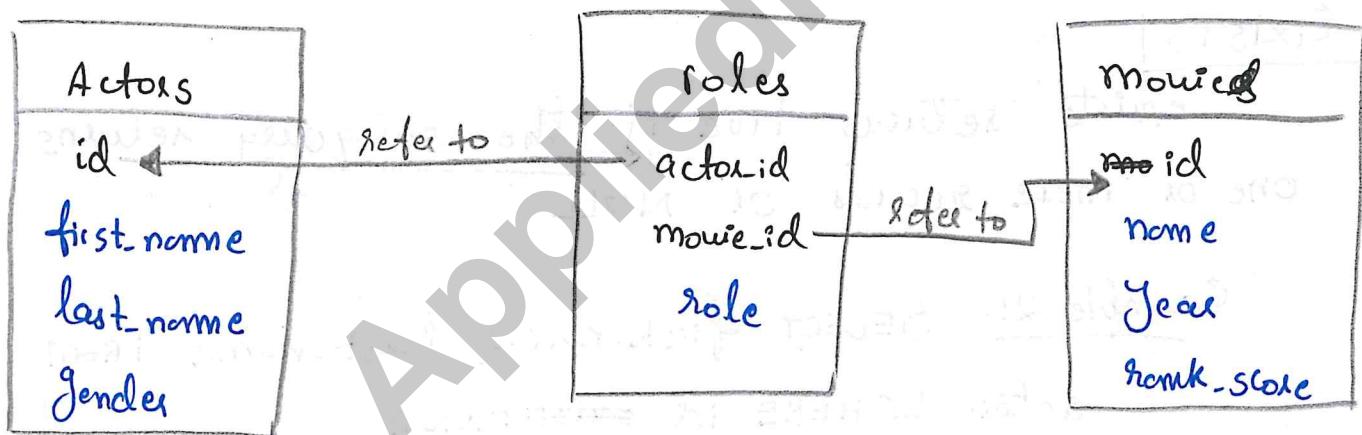
Phone: +91 844-844-0102

SELECT first-name, last-name from actors WHERE
 id IN { SELECT actor-id from roles WHERE
 movie-id IN { SELECT id from movies WHERE
 name = "Schindler's List" } } ; } nested

Keyword to give a select value from a set of list

→ here we are accessing 3 relations in a query using sub-queries. One way to do this is by writing it in the form of subquery while another way to solve this is using "Join's". It will be discussed.

This could be understood as:



It will first execute the innermost query (Q_3)
 i.e "SELECT id FROM movies WHERE name =
 "Schindler's List";"

⇒ It will result in the set of 126 tuples.

Next query which will get executed is Q_2 i.e
 "SELECT actor-id FROM roles WHERE movie-id
 IN (126 tuples);"
 ⇒ It will return set of actor-id's.

At last, the outermost query (Q_1) got executed
 i.e "SELECT first-name, last-name FROM actors
 WHERE id IN (set of all actor id's);"

\Rightarrow It will op all first-name, last-names whose
 the set actor id is in set op of Q_2 :

NOTE

In SQL 'IN' can be thought as belongs to (\in)
 in set theory. While 'NOT IN' can be
 thought as do not belongs to (\notin) in set
 theory.

EXISTS

Exists returns true if the subquery returns
 one or more records or NULL

Example 2: SELECT first-name, last-name FROM
 actors WHERE id ~~exist~~ EXISTS

(SELECT actor-id FROM Roles WHERE movie-id
 EXISTS (SELECT id FROM movies WHERE
 nome = "Schindler's list")));

\Rightarrow If the sub queries return a non-empty
 set say $\{f_1\}, \{f_2\}$, having one or more than
 1 element or return a NULL then what we
 will get whatever we want in the result set

Mail: gatotoco@appliedroots.com

NOTE Similarly it works for NOT EXISTS.

ALL

All operators returns TRUE if all of the subquery values meet the condition.

Example 3: `SELECT * FROM movies WHERE`

`rank-score >= ALL (SELECT MAX(rank-score)`
`FROM movies);`

will return true if ^{all} subquery meet condition.

⇒ It will get us all movies whose rank-score is same as the maximum rank-score.

Here, subquery will return only one value which is `max(rank-score) = {9.9}`. Now outer query will print / select all columns ~~to~~ from table movies where `rank-score >= {9.9}`.

NOTE- Subqueries / inner queries are more easy to read and write when compared to Join's.

ANY

Any operator returns TRUE if any of the subquery set values meets the condition.

* Correlated sub queries \Rightarrow

In SQL, Correlated subquery / synchronized subquery is a subquery that uses values from the outer query. Because the subquery may be evaluated once for each row processed by the outer query, it can be efficient.

Example:

```
SELECT emp-no, name FROM Employees emp
WHERE salary > (SELECT AVG(salary)
                  FROM Employees
                  WHERE department =
                        emp.department);
```

\Rightarrow Here, inner query will return :- it is going to compute average salary of employees and finally will point those emp-no and names of employees whose salary is greater than average salary of Employees within that department. (we know this fact of within that department because we are using alias emp in our inner query).

i.e. per department list all the employees whose salary's are greater than avg salary in that department.

QnB

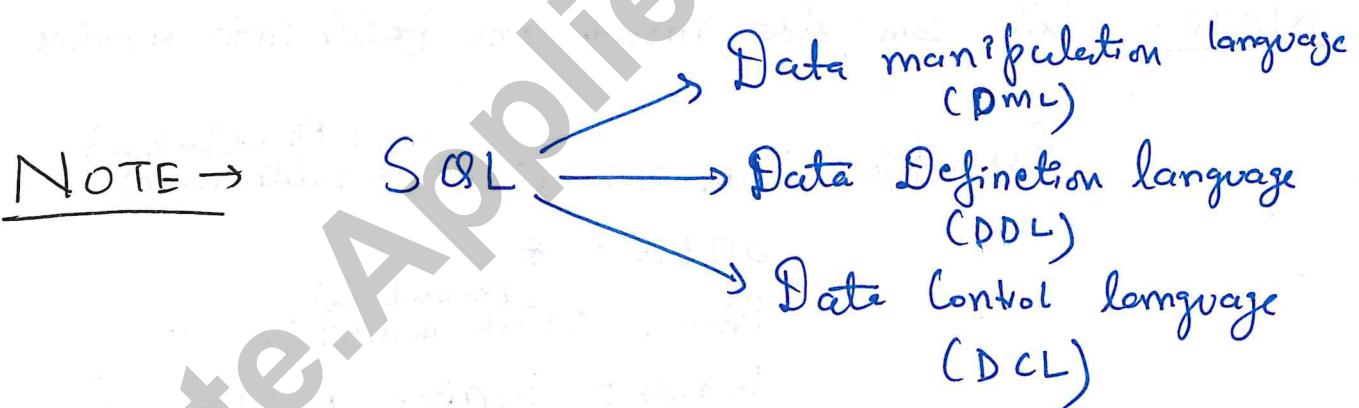
It is called Correlated Subquery because for every employee in outer query we need to run inner subquery again and again.

So, here we using info of outer query in the inner query and this inner query needs to be run for every employee.

NOTE- Correlated queries are computationally expensive.

DATA MANIPULATION LANGUAGE

1.1 Insert:



DML includes SELECT, INSERT, UPDATE, DELETE.

These four commands ~~are~~ belongs to data manipulation language.

INSERT is used to insert data / row / tuples in a relation.

Syntax- `INSERT INTO <table-name> (<attribute-names> VALUES (<values-as-per-attribute-sequence>);`

Example:

`INSERT INTO movies (id, name, year, rank-score)
VALUES`

inserting multiple rows in relation Movies

<code>(412321 , 'THOR' , 2011, 7),</code>	$\rightarrow 1^{\text{st}}$ row
<code>(412323, 'IRON MAN 2' , 2010, 7);</code>	$\rightarrow 2^{\text{nd}}$ row

NOTE - Since 'id' is the primary key in table "movies" So, we cannot insert ~~two~~ rows with same 'id' value.

NOTE - We can also insert one table into another

Example

`INSERT INTO <table-name>`

`SELECT *`

`FROM <table-name>`

`WHERE NAME IN ('John Doe',
'Peter Doe')`

\Rightarrow It will copy all rows from Phone-book to phone-book-2 whose names are either ' John Doe' or ' Peter Doe'.

11.2 Update, Delete:

→ Update:

Syntax: UPDATE <Table-name> SET
 $Col_1 = val_1, Col_2 = val_2 \text{ WHERE}$
 $<\text{Condition}>;$

Example -

UPDATE movies SET rank-score = 9
 $\text{WHERE id} = 412321;$

→ It will update the rank-score value of
 the tuple whose id = 412321.

We can also update multiple column values by
 separating them using commas or using subquery.

→ Delete:

Syntax: DELETE FROM <table-name>
 $\text{WHERE } <\text{Condition}>;$

Example:-

DELETE FROM movies WHERE
 $id = 412321;$

→ It will delete the row which is having
 $id = 412321$

NOTE - If we want to delete all rows from table
 we use a command called TRUNCATE. It is not

DMQL it is a DDL. Its Syntax is: TRUNCATE TABLE
 $<\text{tablename}>;$

NOTE-ED We can also delete all rows from a table using 'DELETE' command. This could happen if we remove WHERE clause.

Syntax - DELETE FROM <table-name>;

DATA DEFINITION LANGUAGE (DDL)

12.1 Create Table:

- * DDL helps us to define, alter, drop a table.
- * To create a new table in database is created as follows:

Syntax: CREATE TABLE <table-name>(<attribute-name> <~~datatype~~> <Specification>);

Example

CREATE TABLE language (

id INT PRIMARY

long VARCHAR(50) NOT NULL);

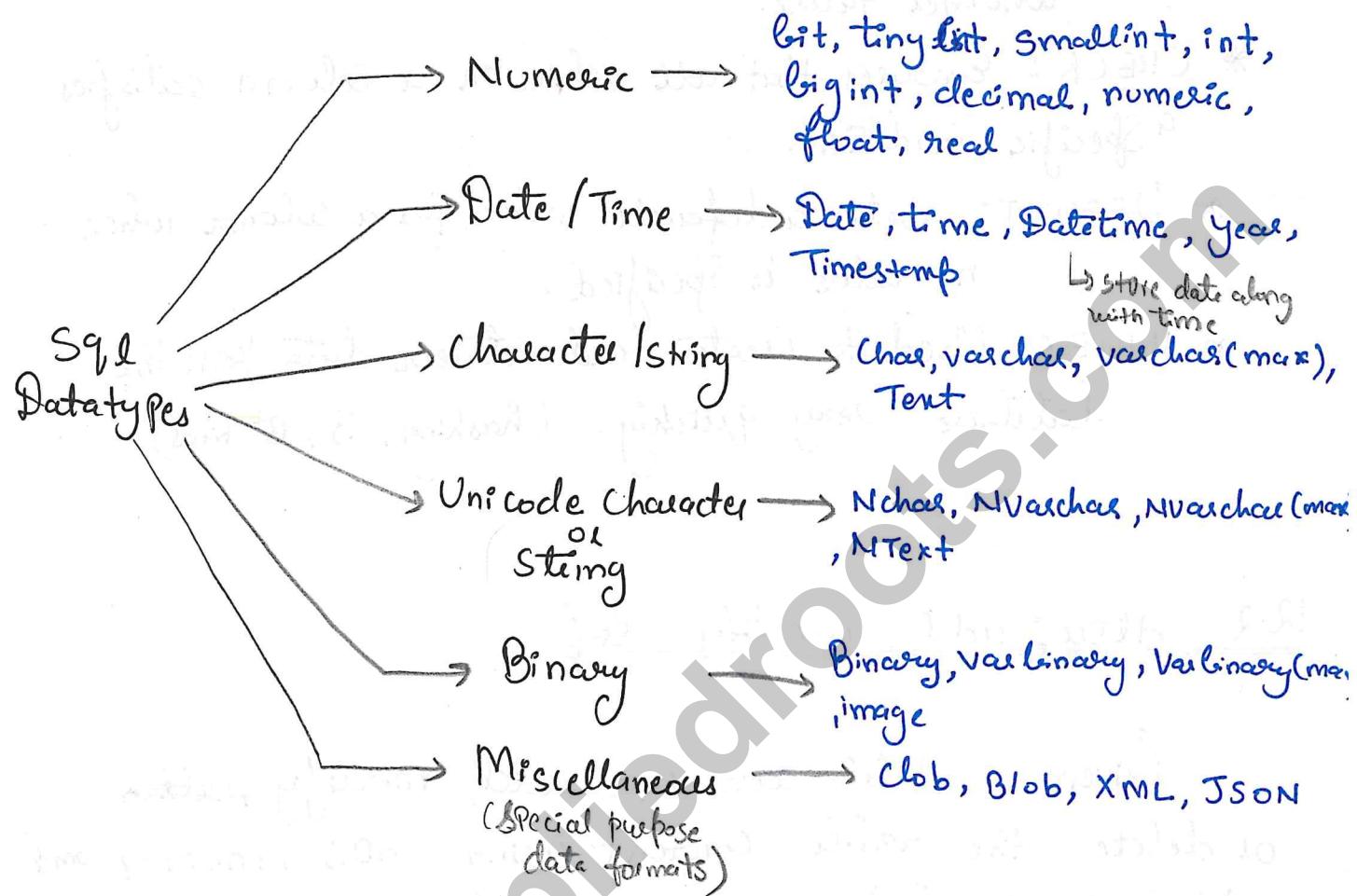
on
Properties
of
Constraint

→ it is integer
type

→ it is a
character
type

⇒ If we try to insert any duplicate id or long as NULL, it will throw an error.

→ SQL Data Types:



Note that, not every database support all data types, Some may support some data types while other may support other data types.

→ Constraints / Properties:

* NOT NULL: Ensures that a column cannot have a NULL value.

* UNIQUE: Ensures that all values in a column are different.

* PRIMARY KEY: A combination of NOT NULL and UNIQUE. It uniquely identifies a row/tuple in a table.

* FOREIGN KEY: Uniquely identifies a row/record in another table.

* CHECK: Ensures that all values in a column satisfies a specific condition.

* DEFAULT: Sets a default value for a column when no value is specified.

* INDEX: Used to create and retrieve data from the database very quickly. (hashing, B, B+ trees)

12.2 Alter; add, modify, drop:

Given a table we can add, modify ~~and~~ or delete the table using ADD, MODIFY and DROP commands in DDL

→ ↕ ADD : It helps us to add new columns or constraints in DB

Syntax: ALTER TABLE <table_name>
ADD <column_name> <data-type>
<constraints>;

Example ALTER TABLE language ADD
Country VARCHAR(50);

→ ↕ MODIFY: It helps to modify column name, datatype or constraints.

Syntax: ALTER TABLE <table_name> MODIFY
<Column_name> <data-type> <constraints>;

Example- ALTER TABLE language ~~MODIFY~~ ^{language: +91 841 841-0102} country
VARCHAR (60);

→ Drop: It will ^(delete) drop the mentioned column.

Example ALTER TABLE language ~~language~~ ^{Drop} country;

→ It will remove whole column "country".

12.3 Drop Table, Truncate, Delete:

→ Drop table: It will remove both the table and all
of the data permanently.

Syntax: DROP TABLE <table name>;

Example: DROP TABLE language;

It will wipe off whole table, while delete
command is use to delete some or all rows of table.
It will never delete the table structure.

NOTE- We can also add one modification to it

DROP TABLE <table name> IF EXISTS;

→ It will drop the table (whole table) if it
exists in database. So that we don't get any
error if we don't know whether it exist or not.

 ~~Applied Roots~~ Truncate: It will remove the content of table but not table structure.

Syntax - TRUNCATE TABLE <table-name>;



DELETE FROM <table-name>

DATA CONTROL LANGUAGE (DCL)

13.1 Grant, Revoke:

In Database, different people have different access to different parts of data because same database is being accessed by multiple people in an organization. Hence, a s/w developer, data scientist, admins etc all have different access permission of the database.

may have permission to add, delete, modify etc

can only access the data (read). (cannot modify it)

have permission to create, insert relation in a database.

DCL helps us to control who can access what. It is not implemented by all databases.

$\rightarrow \leftarrow$ GRANT:

It allows specified users to perform specified tasks.

Example:- GRANT ALL ON db1.* To "jeffrey @'localhost';

↓
all
stands
for all
permissions

db1 is
the name
of database
and * stands
for all tables
of database.

Example 2:- GRANT 'role 1', 'role 2' To 'user1@'localhost',
'user2 @'localhost';

here these are different roles
like admin etc which is
given to user 1 and user 2.

Example 3: GRANT SELECT ON world.* To 'role3';

it has given
Select Permission
on database named
world over all
of its tables / relations

$\rightarrow \leftarrow$ REVOKE:

It cancels previously granted or denied permissions.

Example 1: REVOKE INSERT ON *.* FROM
"jeffrey@'localhost';

here we are revoking
insert permission on all tables
of all database and
all database and
all tables of each database

Example 2: REVOKE SELECT ON world.* FROM 'role 3';

⇒ Revoking Select permission from role 3 on database named world over all of its tables.

NOTE- Best place to learn sql is "reference manual" of any database we are using.

SQL : EXAMPLES AND SOLVED PROBLEMS

- Q2) A relational database contains two tables student and performance as shown below:

Student:

Role no	StudentName
1	Amit
2	Priya
3	Rohan
4	Smita
5	Vinit

Performance:

	Roll_no	Subject Code	Marks
1	1	A	86
2	1	B	95
3	1	C	90
4	2	A	89
5	2	C	92
6	3	C	80

The primary key of the student table is Roll_no. For the performance table, the columns Roll_no and Subject_code together form the primary key. Consider the SQL query given below:

SELECT S.student_name, sum(P.marks)

From Student S, Performance P

WHERE P.marks > 84

Group By S.student_name;

The number of rows returned by the above SQL query is _____

is _____

⇒

Given the query:

SELECT S.student_name, sum(P.marks)

executed 1st From student S, Performance P } → This indicates cross product of S and P

executed 2nd WHERE P.marks > 84 } except 6th row "Performance" is satisfied.

executed 3rd Group By S.student_name; } → this will generate 5 groups with 5 distinct student name

Group by will create 5 groups for 5 distinct student name hence, no. of rows / tuples = 4

(Q2) Consider the following two tables + and four queries in SQL.

Book (isbn, bname), Stock (isbn, copies)

Query 1:

```
SELECT B.isbn, S.copies  
FROM Book B INNER JOIN Stocks  
ON B.isbn = S.isbn;
```

Query 2:

```
SELECT B.isbn, S.copies  
FROM Book B LEFT OUTER JOIN Stocks  
ON B.isbn = S.isbn;
```

Query 3:

```
SELECT B.isbn, S.copies  
FROM Book B RIGHT OUTER JOIN Stocks  
ON B.isbn = S.isbn;
```

Query 4:

```
SELECT B.isbn, S.copies  
FROM Book B FULL OUTER JOIN Stocks  
ON B.isbn = S.isbn;
```

Which one of the queries above is certain to have an output that is superset of the o/p's of the other three queries?

- (A) Query 1
- (B) Query 2
- (C) Query 3
- (D) Query 4

⇒ As we know tuples we get in full outer join is
 Union of $B \bowtie S \cup B \bowtie S \cup B \bowtie S$

Hence, Collect optim is the query with FULL OUTER JOIN
 ie Query 4 (D).

(Q43) Consider a database that has the relational Schema Emp (EmpID, EmpName, DeptName). An instance of the schema Emp and a SQL query on it are given below:

Emp:

EmpID	EmpName	DeptName
1	X4A	AA
2	X4B	AA
3	X4C	AA
4	X4D	AA
5	X4E	AB
6	X4F	AB
7	X4G	AB
8	X4H	AC
9	X4I	AC
10	X4J	AC
11	X4K	AD
12	X4L	AD
13	X4M	AE

Query —

```

SELECT Avg(EC.Num)
FROM EC
WHERE (DeptName, Num) IN
  (SELECT DeptName,
          COUNT(EmpID) AS
           EC(DeptName, Num)
  FROM Emp
  GROUP BY DeptName);
  
```

The output of executing the SQL query is —

Here the Nested query is:

```
SELECT AVG(Ec.Num)
```

```
FROM EC
```

```
WHERE (DeptName) IN (
```

```
    SELECT DeptName,  
          COUNT(EmpID) AS  
          EC(DeptName, Num)  
    FROM Emp  
   GROUP BY DeptName);
```

This will be executed first

This is the alias
for the new
table created
by the inner query

This table EC has
DeptName and Num

↳ Count(EmpId)

The inner query will be executed first. It
will group by the based on DeptName i.e.
AA, AB, AC, AD, AE and Count(EmpId) will
give no. of employees in each department. Hence output
of inner query will be:

Ec:

DeptName	Num
AA	4
AB	3
AC	3
AD	2
AE	1

Now the outer query is asking for Avg on
Nums from resultant relation EC.

$$\text{ie } \frac{4+3+3+2+1}{5} = \frac{13}{5} = 2.6$$

- (Q4) Consider the following database table named
top_scores:

Top_scores:

Player	Country	goals
Klose	Germany	16
Ronaldo	Brazil	15
G Muller	Germany	14
Fontaine	France	13
Pele	Brazil	12
Klinsmann	Germany	11
Kocsis	Hungary	11
Batistuta	Argentina	10
Cubillas	Peru	10
Lato	Poland	10
Hinckes	England	10
T Muller	Germany	10
Rahn	Germany	10

Consider the following query:

SELECT ta.player FROM top-scores AS ta
 WHERE ta.goals > ALL (SELECT tb.goals
 FROM top-scores AS tb
 WHERE tb.country = 'Spain')
 AND ta.goals > ANY (SELECT tc.goals
 FROM top-scores AS tc
 WHERE tc.country = 'Germany'));

The number of tuples returned by the above SQL query is _____

⇒ Here we have 2 sub queries/nested query's.

~~First we will~~ Before start solving the innermost query (③) if we look at the first sub query (②) we can see that it will result empty set because there is no country = "Spain" in our table given. Therefore, it will return an empty set.

Hence, it will always be true because $\text{ALL}(\text{empty set})$

||
TRUE

Since, there is logical AND between ② and ③
 \therefore ② is already TRUE, so now we need to execute ③

The inner most query (③) will return:

The goals made by players whose country = "Germany".

Mail: gatcse@appliedroots.com

O/P →

te. goals
16
14
11
10
10

The condition before logical AND is
 $ta.\text{goals} > \text{Arg}\{16, 14, 11, 10\}$ ie $ta.\text{goals}$
 Should be greater than 10.

Since query ② is ignored as it is TRUE ∴ outer
 query will print Player names whose goals are > 10
 Hence, no. of rows returned = 7

(Q5) Consider the following database table named
water_schemes:

water_Schemes:

Scheme_no	district_name	Capacity
1	Ajmer	20
1	Bikaner	20
2	Bikaner	10
3	Bikaner	20
1	Churu	10
2	Churu	20
		10

The number of tuples returned by the following SQL query is:

With total (name, capacity) as
 ① { Select distinct_name , sum (capacity)
 { From water_schemes
 { group by distinct_name
 { with total_avg (capacity) as
 { Select avg(capacity) as
 { Q2 } from total
 { Select name
 { Q3 } from total, total_avg
 { Where total.capacity >= total_avg.capacity

- (A) 1
- (B) 2
- (C) 3
- (D) 4

→ With is a keyword introduced by oracle. Never mysql version also uses "with".

With creates a temporary relation / table with the result returned by a query and using the given table and attributes name.

Here, Q1 will return \Rightarrow Total:

This is a temporary table which is accessed by the query coming part of further.

name	Capacity
Ajmer	20
Bikaner	40
Churu	30
Dausa	

Sum of Capacity of each city

NOTE - Using "with" we can access the resultant table as 151
it get saved in memory

Phone: +91 844-844-0102

The output of φ_2 will be: (This query is using table)

Total-aug: Created by φ_2)

New table
Created by
 φ_2

Capacity
25

$$\begin{aligned} \text{avg of } & \\ \frac{20+40+30+10}{4} & \\ = 25 & \end{aligned}$$

The output of φ_3 will be: It will use both total and total-aug to get the result. It is performing Cartesian product of total and total-aug and returns the tuples whose capacity is greater than average capacity 25.

hence, only Likanee and Chuu are having Capacity 30 and 40 respectively which is greater than 25
Therefore, no. of tuples returned = 2,

15.2 Solved problems-2:

(Q6) Select operation in SQL is equivalent to

- (A) The selection operation in relational algebra
- (B) The selection operation in relational algebra, except that select in SQL retains duplicates
- (C) The projection operation in relational algebra
- (D) The projection operation in relational algebra, except that select in SQL retains duplicates.

As we know select operation in SQL selects columns while selection operation in RA selects rows/tuples. While projection in RA selects columns but it removes duplicates but select set in SQL retains duplicates.
Hence, correct option is (D).

(Q 7) Consider the following relations:

Students:

<u>Roll_No</u>	<u>Student_Name</u>
1	Raj
2	Rohit
3	Raj

Performance:

<u>Roll_No</u>	<u>Course</u>	<u>Marks</u>
1	Math	80
1	English	70
2	Math	75
3	English	80
2	Physics	65
3	Math	80

SELECT S.Student_Name, SUM(P.marks)

FROM Student S, Performance P }

WHERE S.Roll_No = P.Roll_No }

GROUP BY S.Student_Name;

The number of rows that will be returned by the SQL query is _____

- (A) 0
- (B) 1
- (C) 2
- (D) 3

Mail: gatcse@appliedroots.com

Since, here we are using group by using Student_name and we

(Q8) Consider the following relation:
Cinema (theater, address, capacity)

Which of the following options will be needed at the end of the SQL query:

SELECT P₁.address

FROM Cinema P₁

Such that it always find the addresses of theaters with maximum capacity?

- (A) WHERE P₁.Capacity \geq ALL (SELECT P₂.Capacity FROM Cinema P₂)
- (B) WHERE P₁.Capacity \geq ANY (SELECT P₂.Capacity FROM Cinema P₂)
- (C) WHERE P₁.Capacity $>$ ALL (SELECT MAX(P₂.Capacity) FROM Cinema P₂)
- (D) WHERE P₁.Capacity $>$ ANY (SELECT MAX(P₂.Capacity) FROM Cinema P₂)

\Rightarrow 'Any' will not work here because even if there is even 1 theater whose capacity is max it will get picked up. So, only options left are (A) and (C). When I'm selecting maximum capacity of cinema hall nothing could be greater than it so (C) will not work because ALL will return only one value (which is maximum value) and \therefore it should be \geq and not $>$
Hence, correct and feasible option is (A)

(Q9) Given the following statements. **Phone: +91 844-844-0102**

S₁: A foreign key declaration can always be replaced by an equivalent check assertion in SQL.

S₂: Given the table R(a, b, c) where a and b together form the primary key, the following is a valid table definition.

CREATE TABLE S (

a Integer,

d Integer,

e Integer,

PRIMARY KEY (d),

FOREIGN KEY (a) References R)

Which one of the following statement is CORRECT?

- (A) S₁ is TRUE and S₂ is FALSE
- (B) Both S₁ and S₂ are TRUE
- (C) S₁ is FALSE and S₂ is TRUE
- (D) Both S₁ and S₂ are FALSE

→ S₁: It is not always delete because we do have Cascade option on delete or update.

S₂: In table S, primary key is d and foreign key a, should point to primary key of R. In R (a,b) is a primary key not 'a' alone is a primary key. So, 'a' in table S should refer to (a,b) and not only to 'a'. Hence, it is false

Mail: gatcse@appliedroots.com
Therefore, correct option is (D) Both are false.

(Q10) Given the following Schema:

employees (emp_id, first_name, last_name, hire_date, dept_id, salary)

department (dept_id, dept_name, manager_id, location_id)

You want to display the last names and hire dates of all latest hires in their respective departments in the location ID 1700. You issue the following query:

SQL > SELECT last_name, hire_date

```

    FROM employees
    WHERE (dept_id, hire_date) IN
        (SELECT dept_id, MAX(hire_date)
         FROM employees JOIN departments USING(dept_id)
         WHERE location_id = 1700
         GROUP BY dept_id);
  
```

{ inner query }

What is the outcome?

- (A) It executes but does not give the correct result.
- (B) It executes and gives the correct result.
- (C) It generates an error because of pairwise comparison.
- (D) It generates an error because of the GROUP BY clause cannot be used with table joins in a subquery.

⇒ Since there is no logical and syntactical error in query therefore, it will execute and o/p generated by the query is same as what given in english statement.

Hence it will execute and will give correct answer. Therefore, (B) is the correct answer.

15.3 Solved Problem - 3 :

(Q 11) SQL allows tuples in relations, and correspondingly defines the multiplicity of tuples in the result of joins. Which of the following queries always gives the same answer as the nested query shown below:

SELECT * FROM R WHERE A IN (SELECT S.a FROM S)

- (A) SELECT R.* FROM R, S WHERE R.a = S.a
- (B) SELECT DISTINCT R.* FROM R, S WHERE R.a = S.a
- (C) SELECT R.* FROM R, (SELECT DISTINCT A FROM S) AS S1 WHERE R.a = S1.a
- (D) SELECT R.* FROM R, S WHERE R.a = S.a AND IS UNIQUE R.



Let us consider a sample relation instance of R and S:

R :

a
1
1
2
2
3

S :

a
1
1
2
2

So, the o/p of original query will be:

a
1
1
2
2

Option A will return:

q
1
1
1
1
2

Phone: +91 844-844-0102

Option B will return:

q
1
2

Option C will return:

q
1
1
2
2

Q4 will do join

a
1
2

R
a
1
2
3

Hence, correct option is (C)

(Q 12) Consider the following relational schema:
Employee (empid, empName, empDept)
Customer (custid, custName, salesRepId, rating)
SalesRepId is a foreign key referring to empid of
the employee relation. Assume that each employee
makes a sale to at least one customer. What does the
following query return?

SELECT empName

From Employee E

WHERE NOT EXISTS (SELECT custid
From Customer C

WHERE C.SalesRepid = E.empid
ANDS C.rating <> 'GOOD');

- (A) Names of all employees with at least one of their customers having a 'GOOD' rating.
- (B) Names of all employees with at most one of their customers having a 'GOOD' rating.
- (C) Names of all the employees with none of their customers having a 'GOOD' rating.
- (D) Names of all the employees with all their customers having a 'GOOD' rating.

\Rightarrow The assumption made in question "Each employee makes a sale to at least one customer" means that in SalesRepid, every empid exists at least once.

The inner query will give : custid for SalesRepid = empid and rating is not equal to "good".
 Outside the inner query, we have NOT EXISTS . So, its kind of like negation over negation (as in inner query we have rating != "Good"). So, any employee who have even 1 customer whose rating is not equal to good will not be there in the result set.
 So, option A, B and C are wrong. Hence, D is the most appropriate option.

(Q 13) Consider the following tables 9A, 84B-044C, and 0102

Table A :

Id	Name	Age
12	Arun	60
15	Shreya	24
99	Rohit	11

Table B :

Id	Name	Age
15	Shreya	24
25	Hari	40
98	Rohit	20
99	Rohit	11

Table C :

Id	Phone	Area
10	2200	02
99	2100	01

How many tuples does the result of the following SQL query contains?

```
SELECT A.Id  
FROM A  
WHERE A.age > ALL (SELECT B.age  
                      FROM B  
                     WHERE B.name = 'arun')
```

- (A) 4
- (B) 3
- (C) 0
- (D) 1

⇒ The inner sub query will return empty set and we know $\text{ALL}(\text{empty set}) = \text{True}$
Hence, $A.\text{age} > \text{ALL}()$ is true for all tuples.
and it will return $A.\text{id}$ for all $A.\text{age} > \text{all}$. Therefore,
it will return 3 tuples.

(Q14) Database table by name loan_Records is given below:

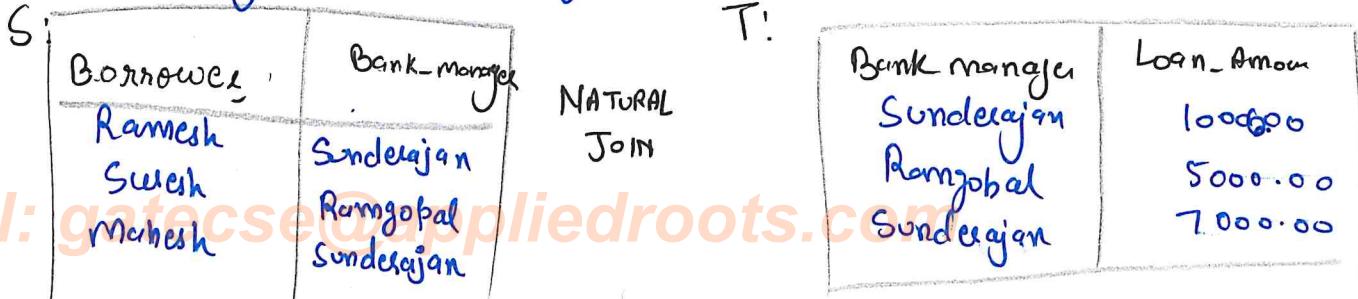
Borrower	Bank_manager	loan_Amount
Ramesh	Sunderajan	10000.00
Suresh	Ramgopal	5000.00
Mahesh	Sunderajan	7000.00

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

```
SELECT COUNT(*)
FROM ((SELECT Borrower, Bank_Manager
      FROM loan_Records) AS S
      NATURAL JOIN (SELECT Bank_manager, loan_Amount
                    FROM loan_Records) AS T);
```

- (A) 3
- (B) 9
- (C) 5
- (D) 6

⇒ According to the query we will have

S: 

NATURAL JOIN

Borrower	Bank_manager
Ramesh	Sunderajan
Suresh	Ramgopal
Mahesh	Sunderajan

Bank_manager	loan_Amount
Sunderajan	10000.00
Ramgopal	5000.00
Sunderajan	7000.00

It will result in distinct 5 tuples :

Borrower	Bank_manager	loan_Amount
Ramesh	Sunderajan	10k
Ramesh	Sunderajan	7k
Buvesh	Rangopal	5k
Mahesh	Sunderajan	10k
Mahesh	Sunderajan	7k

5
distinct
tuples

Hence, output will be 5 tuples (as it's Count(*))

FUNCTIONAL DEPENDENCIES

Phone: +91 844-844-0102

AND

DECOMPOSITION

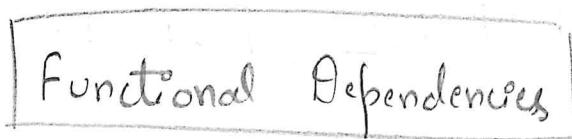
16.1 Designing tables in DB:

Here, database corresponds to the Relational database.

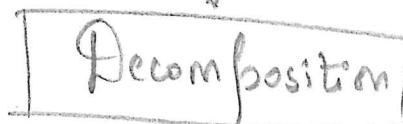
The objectives we carry while designing a database is:

- * From ER diagrams \rightarrow Tables, our objective is to minimize the number of tables.
- * Another objective while designing a db is minimizing storage or disk space while not discarding any data/information. (This could be done by avoiding duplicate data, it is one of the way to minimize disk storage)
- * One of the main objective while designing database is fast retrieval of data. In current scenario the storage cost has fallen drastically while demand for speedy retrieval of data.

16.2 Functional Dependencies Phone: +91 844-844-0102



lays foundation
for



(breaking data
into tables)



to perform
decomposition there
are various form called
normalization forms required

Consider, a real world example to understand
the concept of functional dependencies

Customer (id, name, zip)

There are certain properties of each relations which
are not given till now (implicitly). like:

⇒ Given an 'id' we can uniquely determine name
of a customer. It is represented as:

$$id \rightarrow name$$

⇒ Given an 'id' we can uniquely determine
zip code of a customer.

$$id \rightarrow zip$$

⇒ Given an 'id' we can uniquely determine
both name and zip code of a customer. i.e

$$id \rightarrow name, zip$$

Note Here, ' \rightarrow ' is a notation which means LHS functionally determines RHS.

It means, 'id' functionally determines norm and zip.

In general we can give it as:

Given relation, $R(A_1, A_2 \dots A_n)$

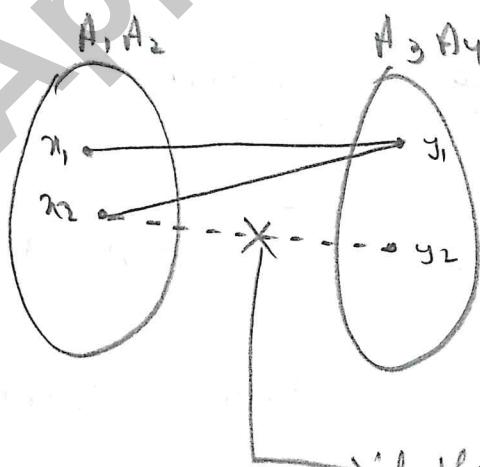
Let

$$\{A_1, A_2\} \rightarrow \{A_3, A_4\}$$

Here, for given values of A_1 and A_2 we can functionally determine A_3 and A_4 .
(Uniquely determine)

Formally, Functional dependency is defined as a constraint between sets of attributes in a relation.

Another way to understand FD is set theoretic way:



If this line exists than for give value x_2 we cannot uniquely determine A_3, A_4 and hence, this FD will fail. So, in order to hold this FD this dotted line cannot exist.

NOTE

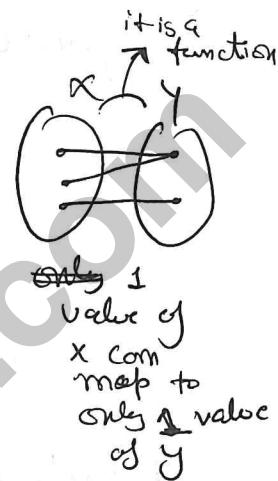
logically: $x \rightarrow y$ ~~Photo: (t₁.x = t₂.x) \rightarrow t₁.y = t₂.y~~ $\Rightarrow 102$

APPLIED
ROOTS

$(t_1.y = t_2.y))$ where t_1 and t_2 are two different tuples.

Mathematically: $T_{x,y}(R)$ is a function

$$x \rightarrow y$$

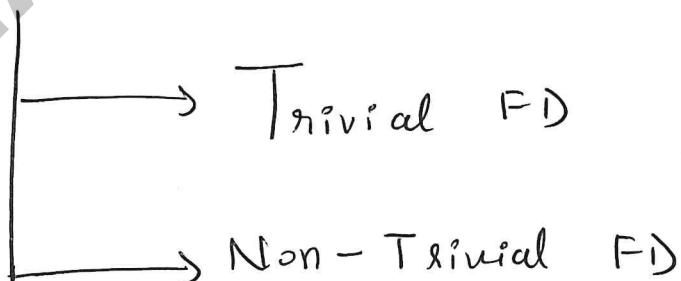


In practice, FD's help us to encode business logic.

For ex- Given a name we cannot uniquely determine an id of a customer

name \rightarrow id (There could be 2 same name with diff id)

Types of Functional dependencies:



(I) TRIVIAL FD's:

For given relation Customer (id, name, zip), the trivial FD's will be:

$$id \rightarrow id$$

$$id, name \rightarrow id$$

$$id, name \rightarrow name$$

} all are trivial FD's

Mail: gatecse@appliedroots.com

$x \rightarrow y$

An FD, $x \rightarrow y$, is called trivial FD in relation R , if X is a superset of Y ($X \supseteq Y$). In previous 3 examples the L.H.S is the superset of R.H.S. Hence, they are trivial.

(II) Non-Trivial FD's:

Given a relation R , an FD $X \rightarrow Y$ is said to be non-trivial if $X \subsetneq Y$ (X is not a superset of Y) ($X \supsetneq Y$).

For the customer relation, FD $id \rightarrow name$ is a non-trivial functional dependency.

\Rightarrow Inference Rules:

- * Based on real world problem / system that we are modelling, we derive / get some FD's (Based on ground reality, or system we get some FD's). (Properties of system)

For example: Multiple people can have same name,
 (this is ~~#~~ one of the property of the system)
 hence we cannot have $name \rightarrow id$

- * Based on the predefined properties of the system and existing FD's we can define some rules to obtain new FD's. These are as follows:

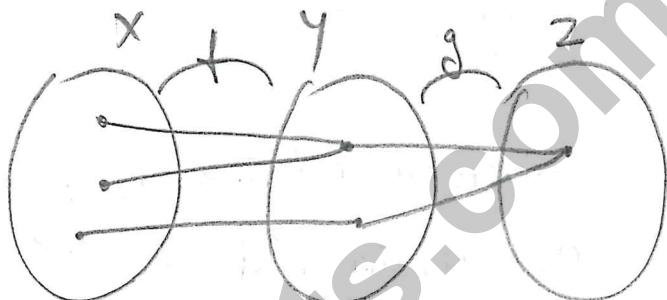
Assume, there is Relation R with n_1 ... as its attributes

$R(x_1, y_2)$

Phone: +91 844-844-0102

$\rightarrow \leftarrow$ Reflexivity \Rightarrow if $x \geq y$ (x is superset of y) then
 $x \rightarrow y$.

$\rightarrow \leftarrow$ Transitivity \Rightarrow if $x \rightarrow y$ and $y \rightarrow z$ then
 $x \rightarrow z$



fog \Leftrightarrow got

$\rightarrow \leftarrow$ Augmentation \Rightarrow if $x \rightarrow y$
then $x_3 \rightarrow y_3$ ($z \rightarrow z$)

$\rightarrow \leftarrow$ Decomposition \Rightarrow if $x \rightarrow y_3 \Rightarrow x \rightarrow y$
and
 $x \rightarrow z$

$\rightarrow \leftarrow$ Union \Rightarrow if $x \rightarrow y$ and $x \rightarrow z \Rightarrow x \rightarrow y_3$

$\rightarrow \leftarrow$ Composition \Rightarrow if $x \rightarrow y$ and $z \rightarrow w \Rightarrow x_3 \rightarrow yw$
(here x, y, z, w all are attributes of
relation R)

$\rightarrow \leftarrow$ Pseudo-transitivity \Rightarrow if $x \rightarrow y$ and $wy \rightarrow z \Rightarrow wz \rightarrow wy$
(Augmentation +
transitivity)

$wx \rightarrow z$

Mail: gatecse@appliedroots.com

NOTE- A common mistake committed by people is:
 Given $xy \rightarrow z$, they decompose LHS, which is wrong. We always decompose RHS.
 Hence $xy \rightarrow z \Rightarrow x \rightarrow z$ and $y \rightarrow z$

Example:

R:

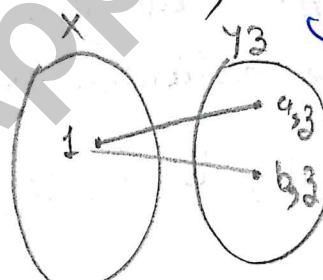
X	Y	Z
1	a	a
1	b	a
2	a	a
2	b	a

here the non-trivial dependency which will hold are:

$$x \rightarrow z$$

(given X we can uniquely determine Z)

$$\Rightarrow x \rightarrow y \text{ and } x \rightarrow yz$$



These does not hold because given X we cannot uniquely determine Y or Yz.

$$\Rightarrow y \rightarrow z \quad (\text{it will hold} \text{ (valid)})$$

$$\Rightarrow y \rightarrow x \text{ and } y \rightarrow xz \quad (\text{it does not hold} \text{ (not valid)})$$

$$\Rightarrow xy \rightarrow z \quad (\text{it will hold} \text{ (valid)})$$

$$yz \rightarrow x \quad \left. \begin{array}{l} (\text{it will not hold}) \\ \text{and} \\ x \rightarrow z \end{array} \right\} (\text{not valid})$$

Other, non-trivial FD's are $X_4 \rightarrow Y_3$, $Y_2 \rightarrow Z_2$.

16.3 Attribute closure, Keys and Solved problems:

⇒ Attribute closure :-

It is a set of attributes that can be functionally determined using attributes in X . (R is a relation with X as a set of attributes). It is represented as ' X^+ ' (here, X can be single att. or multiple attribute (set))

Example 1: Given relation $R(A, B, C, D)$ and

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

Set of FD's

⇒ now A^+ will give all the attributes we can determine or get using A

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

↳ (Because it is, trivial)

Similarly, $B^+ = \{ B, C, D \}$ (using B we can determine B, C and D)
 $C^+ = \{ C, D \}$ (using C we can determine D)

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

Example 2:

Given relation $R(A, B, C, D, E, F, G)$ and
 $FD = \{AB \rightarrow CD, AF \rightarrow D, DE \rightarrow F, G \rightarrow G, F \rightarrow E, G \rightarrow A\}$

Now if we want to find closure of $CF (\{CF\}^+)$
it will be given as:

$$\{CF\}^+ = \{C, F, E, \overset{\text{trivial}}{\underset{\substack{\uparrow \\ \text{Using } C}}{G}}, A, D\}$$

\downarrow Using F we can determine E ($F \rightarrow E$) \downarrow \because we got G and $G \rightarrow A$ \downarrow \because we got A and $AF \rightarrow D$

\uparrow \because A is there and hence $AB \rightarrow CD$

Similarly $\{BG\}^+ = \{B, \overset{\text{trivial}}{\underset{\substack{\uparrow \\ \text{Using } G \\ G \rightarrow A}}{G}}, A, C, D\}$

$$\{AF\}^+ = \{A, \overset{\text{trivial}}{\underset{\substack{\downarrow \\ \text{using } A \text{ and } F \\ AF \rightarrow D}}{F}}, D, E\}$$

\downarrow $\because F \rightarrow E$

M-Sub

* Application of Attribute closure concept is to determine super keys and candidate keys. If we have a relation and FD for a system then we can determine all of its candidate key using this concept.

If $X^+ = \text{all attributes in } R \Rightarrow$ uniquely determine a tuple/row
then X is a superkey.

Mail: gatcse@appliedroots.com A minimal superkey is a candidate key. Minimal means, if no subset or strict subset of a superkey is a key,

then, that superkey is a candidate key.

That is, if X is a superkey $\Rightarrow X^+ = \text{all attributes of } R$

then, X is a candidate key iff $\forall Y \subset X$
 $Y^+ \neq \text{all attributes of } R$.

Hence, we can pick any of the candidate key to be the primary key.

Example 3:

Given relation $R(ABCDE)$ and

$FD = \{ A \rightarrow BC, D \rightarrow E \}$. What are the ~~are~~ superkeys (SK) and candidate keys (CK)

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

↑
Superkey.

Now this superkey to be the candidate key, no subset of it i.e. A and D should uniquely determine all attributes. Therefore we will check A^+ and D^+

$$\left. \begin{array}{l} A^+ = \{ A, B, C \} \\ D^+ = \{ D, E \} \end{array} \right\} \text{not superkeys}$$

Hence, no subset of $\{AD\}$ is a superkey. Therefore $\{AD\}$ is the candidate key.

Example 4:

Given relation $R(A B C D E)$ and

$FD = \{ A B \rightarrow C, C \rightarrow D, B \rightarrow E \}$. what is the CK?

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

$$\begin{aligned} \{A\}^+ &= \{A\} \\ \{B\}^+ &= \{B, E\} \end{aligned} \quad \text{Not a super key.}$$

Hence, AB is a candidate key.

Ques

NOTE- Attributes which are part of atleast one candidate key are called as prime attribute. In the above example 4, A and B are prime attributes while rest of the attributes are non-prime attributes.

The prime attribute concept is important because:

lets understand it using an example. Given relation $R(A, B, C, D, E)$ and FD's $= \{ A B \rightarrow C, C \rightarrow D, B \rightarrow E, E \rightarrow B \}$

\Rightarrow here CK will be as

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

$$\text{while } A^+ = \{A\} \text{ and } B^+ = \{B, E\}$$

So here prime attributes are $\{A, B\}$

In FD's we can see there is functional dependency $E \rightarrow B$ so whenever we have any FD such that

$X \xrightarrow{\text{prime attribute}}$
then, we may have more candidate keys.

So, if we replace that prime attribute Phone_1 with $X44-0102$
then it will also become a candidate key.

i.e. $AE^+ = \{A, B, C, D, E\}$ is a candidate Key

as $A^+ = \{A\}$ and $E^+ = \{E, B\}$. and prime attributes
are $\{A, B, E\}$ now

Hence, this will have 2 candidate keys AB and AE .

$$CK = \{AB, AE\}$$

So, concept of prime attribute helps us to determine all possible candidate keys.

Example 5:

$R(ABCDE)$

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

$CK?$



$AB^+ = \{A, B, C, D, E\}$ } This is not a CK because subset of it can uniquely determine all attributes.

$A^+ = \{A\}$

$B^+ = \{B, A, C, D, E\}$ } ?? it can determine all attributes.
~~it = { }~~ ∴ it is a CK

So prime attributes are: $\{B\}$

But, here we have $FD = C \rightarrow B$ (i.e. $x \rightarrow$ prime attribute)

hence C is also a CK

Hence, Prime attributes are = $\{B, C\}$

NOTE- C is being determined by both AB ($AB \rightarrow C$)
but AB could not be a CK because B itself is a CK. $\therefore AB$ is a SK.

Example 6:

R(A B C D E)

Phone: +91 844-844-0102

FD = { A → BCDE, BC → ADE, D → E }

CK?

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

Here A is also being determined by BC
 $BC \rightarrow ADE$

∴ we will check for BC also

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

∴ $B^+ = B$, $C^+ = C$ ∴ BC is a ck

Hence, prime attributes are : {A, B, C}
 candidate keys are : {A, BC}

Example 7:

R(A B C D)

FD = { A → B, B → C, C → D, D → A }

CK?

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

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

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

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

Hence, all A, B, C, D are ck.

This type of FD is called cyclic FD.

Example 8:

$R(A B C D E)$

Phone: +91 844-844-0102

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

CK?

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

$A B$ is a ck $B \in A^+ = A, B^+ = B$

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

$A F$ is a ck $B \in F^+ = F B$

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

$A E$ is also a ck $B \in E^+ = E, B, F$

Hence CK = {AB, AF, AE, AD, AC}

and prime attributes are = {A, B, E, F, D, C}

Similarly for AD and AC.

Example 9:

$R(A B C D E)$

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

CK?

\Rightarrow Here, E is not part of any FD



E has to be part of every
Consolidate Key as
it could not be determined by
any other attribute

Mail: gatecse@appliedroots.com
 $A B E^+ = \{ A, B, E, C, D \}$ is a ck because
 $A E^+ = A, E$ and $B E^+ = B, E$

Example 10:

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

$FD = \{ AB \rightarrow CDEF, C \rightarrow ADE, D \rightarrow BEF, F \rightarrow DA, BE \rightarrow AF \}$

CK?

\Rightarrow Here, G is not in any FD $\Rightarrow G$ has to be part of every CK.

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

$$AC^+ = \{A, G\}$$

$$BC^+ = \{B, G\}$$

Hence ABG is a CK.

So,

\therefore we have a FD $C \rightarrow ADE$ ($X \rightarrow$ Prime attribute)

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

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

$$BG^+ = \{B, G\}$$

$\therefore CBG$ is ~~not~~ a CK. CG is a CK

we also have FD $D \rightarrow BEF$ ($X \rightarrow$ Prime attribute)

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

$$\text{Similarly } BEG^+ = \{B, E, G, A, C, D, F\}$$

like this we have multiple CK here

$$\{ABG, CG, DG, FG, BEG\}$$

$$\text{Prime attribute} = \{A, B, C, D, E, F, G\}$$

16.4 Functional dependency

Sets ~~and~~: Properties 91-844-0102



Solved Examples:

→ Properties of functional dependency set:

→ Membership:

Let F be a FD set on R , if $x \rightarrow y$ using FD's in F then, $x \rightarrow y$ is a member of F

here, if $x \rightarrow y$ using FD's in F means $x^+ \supseteq y$ then using FD in F .

Example 1 $R(A B C)$ and $F = \{A \rightarrow B, B \rightarrow C\}$

(Q1) Is $A \rightarrow C$ is a member of F ?

⇒ $A^+ = \{A, B, C\} \Rightarrow$ Using FD in set F
we can say
that $A \rightarrow C$
which means $A \rightarrow C$
is member of F .

(Q2) Is $C \rightarrow A$ is a member of F ?

⇒ To prove this true we should be able to get $C \rightarrow A$ using FD's in F .

$C^+ = \{C\} \Rightarrow$ using FD's in F it can give only C and they cannot determine any other attribute.

Hence $C \rightarrow A$ and it is not member of F .

Mail: gatecse@appliedroots.com

Example 2:

Given $R(A, B, C, D, E)$ and $FD = \{AB \rightarrow C, A \rightarrow D, C \rightarrow E\}$

Is $AC \rightarrow E$ is a member of FD set F?

$\Rightarrow A^+ = \{A, C, D, E\} \Rightarrow \therefore AC \text{ can determine } E$
 $\therefore AC \rightarrow E \text{ holds and it is a member of } F.$

Example 3:

Given two conditions $A \supseteq B$ and $B \rightarrow c$
 does $A \rightarrow c$ holds?

\Rightarrow Here we can construct our FD set. Since,
 A is superset of B

$\therefore A \rightarrow B$

and $B \rightarrow c$ is already given

$\therefore FD = \{A \rightarrow B, B \rightarrow c\}$

$\therefore A^+ = \{A, B, C\} \Rightarrow$ hence c could be determined by A
 and $\therefore A \rightarrow c$ is member of FD set F.

Example 4:

$R(A, B, C, D)$ and $FD = \{AB \rightarrow C, BC \rightarrow D\}$

Does $A \rightarrow c$ is a member of FD set F?

$\Rightarrow A^+ = \{A\} \Rightarrow c$ could not be determined by A
 hence $A \rightarrow c$ does not hold and
 hence, it is not a member of F.



Closure of a FD set :

It is also written as F^+ . It is a set of all FD that can be determined using FD in F.

Example 1:

$R \in A B C \}$ and $F = \{ A \rightarrow B, B \rightarrow C \}$
What is F closure (F^+)?

$\Rightarrow F^+$ will include all possible trivial and non-trivial dependencies.

$$F^+ = \left\{ \begin{array}{c|c|c} A \rightarrow A & AB \rightarrow A & \dots \\ A \rightarrow B & AB \rightarrow B & \dots \\ A \rightarrow C & AB \rightarrow C & \dots \end{array} \right\}$$

It is a very large set

A question which used to get framed from this concept is

(Q) Number of FD's in F^+ ?

\Rightarrow For the same example $R(A B C)$ and $F = \{ A \rightarrow B, B \rightarrow C \}$ we will try to get total number of FD in F^+

* FD's with 0 attributes $\Rightarrow \emptyset \rightarrow \emptyset$

// valid FD
bcz any
dependency to be FD
 $x \rightarrow y$
 $x \geq y$.

This
subset could
be empty.

So, 1 FD is possible
with 0 attributes.

* FD's with 1 attribute
 (by 1 attribute we mean LHS is having a single attribute) in FD

Phone: +91 844-844-0102

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

any subset of it

$$= 2^3 = 8$$

$A \rightarrow B$
 $A \rightarrow C$
 $A \rightarrow BC$

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

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

any subset of it

$$= 2^2 = 4$$

$B \rightarrow B$
 $B \rightarrow C$
 $B \rightarrow BC$

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

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

any subset of it

$$= 2^1 = 2$$

So, total $8 + 4 + 2$ FD's are possible with 1 attribute in LHS

* FD's with 2 attributes $\rightarrow AB^+ = \{A, B, C\}$

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

all subset of it

$$= 2^3 = 8$$

$AB \rightarrow A$
 $AB \rightarrow B$
 $AB \rightarrow C$
 $AB \rightarrow ABC$

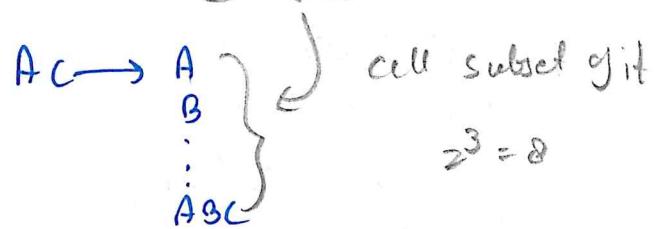
$$BC^+ = \{B, C\}$$

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

all subset of it

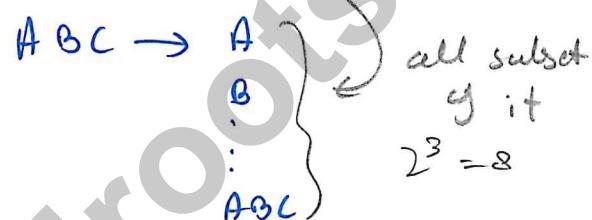
$$= 2^2 = 4$$

$BC \rightarrow B$
 $BC \rightarrow C$



So, total $8+4+8$ FD's possible with 2 attributes in LHS

* FD's with 3 attributes $\Rightarrow A \rightarrow^+ \{A, B, C\}$



So, 8 FD's are possible with 3 attributes in LHS

Hence, total no. of FD's in F^+ given = $8+8+4+8+8+4+2+1$
 $= \boxed{43 \text{ FD's}}$

So, ~~on~~ on seeing F, among 43 FD's we need only 2 FD's to determine rest of 41. Therefore, these 41 FD's are called redundant FD's as we can derive these 41 just using these 2 FD's ($f: A \rightarrow B, B \rightarrow C$).

$\Rightarrow \Leftarrow$ Equality of FD set \Leftrightarrow Phone: +91 844-844-0102



Given two functional dependency set F and G,
equality of FD set means that F and G are
equal iff

(i) $F^+ = G^+$ (But it is very tedious to compute)

OR

(ii) F covers G and G covers F

it means if all FD's in G can be determined using FD's in F and viceversa.

Note
It can be thought as 2 diff engineer come up with 2 FD set based on their understanding. Now to check if both set are equivalent this property is used.

Example :-

R(A B C D)

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

$$G = \{AB \rightarrow C, AB \rightarrow D, C \rightarrow D\}$$

Is $F = G$?

\Rightarrow F covers G?

$$AB_F^+ = \{AB, C, D\} \Rightarrow AB \rightarrow C \text{ and } AB \rightarrow D \text{ in } G$$

(it means AB^+ under F)

$$C_F^+ = \{C, D\} \Rightarrow C \rightarrow D \text{ in } G$$

Hence, F Covers G.

G covers F?

Phone: +91 844-844-0102



$$A^+_{\bar{G}} = \{A, B, C, D\} \Rightarrow \underbrace{A \rightarrow C}_{\text{in } F}, \underbrace{D}_{\text{not in } F}$$

(it means $A \rightarrow C$ under G)

$$B^+_{\bar{G}} = \{B\} \Rightarrow \underbrace{B \rightarrow C}_{\text{in } F}$$

$\because F$ have $B \rightarrow C$ but
G does not hold it

$$C^+_G = \{C, D\} \Rightarrow \underbrace{C \rightarrow D}_{\text{in } F}$$

$\therefore G$ does not cover F.

$\therefore F \neq G$

(For F and G to be equivalent both
F and G should cover each other)

Example 2:

R(A B C D E F H)

$$F = \{A \rightarrow C, AC \rightarrow D, E \rightarrow AD, E \rightarrow H\}$$

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

F \subseteq G ?

$\Rightarrow F$ covers G -

$$A^+_{\bar{F}} = \{A, C, D\} = \underbrace{A \rightarrow CD}_{\text{could be determined using } F}$$

$$E^+_{\bar{F}} = \{E, A, D, H, C\} = E \rightarrow AH$$

G covers F -

$$A^+_{\bar{G}} = \{A, C, D\} \Rightarrow \underbrace{A \rightarrow C}_{\text{in } F}$$

$$AC^+_{\bar{G}} = \{A, C, D\} = \underbrace{AC \rightarrow D}_{\text{in } F}$$

$$E^+_{\bar{G}} = \{E, A, H, C, D\} = \underbrace{E \rightarrow AD}_{\text{and}} \quad \underbrace{E \rightarrow H}_{\text{in } F}$$

Mail: gatetcse@appliedroots.com

| Hence, F \subseteq G |

Hence, G covers F

16.5 Minimal / Canonical Covers of FD's + Solved Problems:

Objective: Eliminating Functional dependencies which are redundant in an FD-set F.

Redundant FD's are those whose removal does not impact F^+ .

Example 1:

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

What is minimal cover of given set F.

\Rightarrow

$$A^+ = \{A, B, C\} \xrightarrow{\text{using } A \rightarrow B}$$

$$B^+ = \{B, C\} \xrightarrow{\text{using } B \rightarrow C}$$

$$\text{So, } F_1 = \boxed{\begin{matrix} A \rightarrow B \\ B \rightarrow C \end{matrix}} \xrightarrow{A \rightarrow C}$$

$$\therefore F^+ = F_1$$

Hence, $A \rightarrow C$ is redundant as we can get C using the remaining FD's.

So this F_1 is called minimal cover of F as it is the smallest set of dependency possible to get all the FD's preserved.

Ques

The need of minimal cover is to minimize the number of FD's. In SQL FD's are enforced using CHECK and CONSTRAINTS. So smaller the no. of FD's lesser the CHECK operation to be performed and faster the database queries will be executed.

NOTE

Phone: +91 844-844-0102

Entaneous variable /attribute :- The variable / attribute which is not responsible is said to be entaneous, if can we remove it and its removal does not bring any change in the closure of the set of FD's.

Ex - $R(ABC)$

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

\Rightarrow here B is entaneous variable

$$\text{minimum cover of } F = \{ A \rightarrow C \}$$

Getting rid of all entaneous attribute is a way to obtain a minimal cover of FD set F.

Key idea is : if we have $X \rightarrow w$ and $XY \rightarrow w$ then we can easily remove Y as we can already determined w using X. So removal of this attribute in a FD does not change the expressive power of F.

$$\text{i.e. } F^+ = F_1^+$$

Task : Given FD sets find minimal covers of F.

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



Step 1: Find FD's that have some LHS and combine them using Union principle.

$\Rightarrow A \rightarrow BC$
 $A \rightarrow B \}$ combine them using union

So, now new

$$F_1 = \{ A \rightarrow BC, B \rightarrow C, AB \rightarrow C \}$$

$$A \rightarrow BC$$

$\} \text{ in place of having both } A \rightarrow BC \text{ and } A \rightarrow B \text{ we will have } A \rightarrow BC$

Mail: jatecse@appliedroots.com

Step 2: Find FD's that have extraneous attribute variables in LHS or RHS

⇒ we have $A \rightarrow c$ and $B \rightarrow c$

∴ B can uniquely determine C than we don't need AB to determine C

So, now new FD set will be

$$F_2 = \{ A \rightarrow Bc, B \rightarrow c \}$$

Step 3: Remove any extraneous attri, but which you can find now.

⇒ we have $A \rightarrow Bc$ and $B \rightarrow c$

∴ Using A we can determine B and C and using B we can determine C

That means ~~A~~ A is not only one who determine C. Therefore, even if we drop c from $A \rightarrow Bc$, it will not ~~be~~ bring any change
Hence, new FD set will be

$$F_3 = \{ A \rightarrow B, B \rightarrow c \}$$

Note

Here, we can notice that in step 3 we have not decreased/ removed any FD, we can see here that we are even minimizing the individual FD's also. So using minimal cover concept we not only ~~decrease~~ remove the redundant FD's we also minimize the individual FD's.

Step 4: Repeat step 1, 2 and 3 till minimal cover is obtained (i.e. F_1 doesn't change)

⇒ The minimal cover obtained in step 3 is the minimal cover of F as we cannot remove any FD further and cannot minimize the individual FD's also.

Hence $F^+ \Leftrightarrow F_3^+$

Example 2: $F = \{ABCD \rightarrow E, E \rightarrow D, A \rightarrow B, AC \rightarrow D\}$

⇒ Step 1: No common/same LHS FD's are there.
⇒ Step 2:

$$\begin{array}{l} A \rightarrow B \\ \cancel{ABCD} \rightarrow E \end{array} \Rightarrow \begin{array}{l} A \rightarrow B \\ ACD \rightarrow E \end{array}$$

We can remove B from $ABCD \rightarrow E$ because given A we can uniquely determine B .

$$F_1 = \{A \rightarrow B, ACD \rightarrow E, E \rightarrow D, AC \rightarrow D\}$$

Step 3:

$$\begin{array}{l} ACD \rightarrow E \\ AC \rightarrow D \end{array} \Rightarrow \begin{array}{l} AC \rightarrow D \\ AC \rightarrow E \end{array}$$

Now these 2 have some LHS. Hence we can combine them.

$$\begin{array}{l} AC \rightarrow E \\ AC \rightarrow D \end{array} \Rightarrow AC \rightarrow DE$$

$$F_2 = \{AC \rightarrow DE, E \rightarrow D, A \rightarrow B\}$$

Step 4:

$$\begin{array}{l} E \rightarrow D \\ AC \rightarrow DE \end{array} \quad \left. \begin{array}{l} E \rightarrow D \\ AC \rightarrow E \end{array} \right\}$$

$$F_3 = \{E \rightarrow D, AC \rightarrow E, A \rightarrow B\}$$

* Step 5: Common mistakes

People try to ~~forget~~ drop c from $AC \rightarrow E$. But it will be wrong

$$\text{ie } F' = \{E \rightarrow D, A \rightarrow E, A \rightarrow B\} \times$$

because E could be determined only using A and c and not only using A.

In this case $F \neq F'$

(minimal cover should not impact F')

So, always verify at the end whether the minimal covers ~~both~~ F and vice versa

Hence, the minimal cover of F is F_3

$$\{E \rightarrow D, AC \rightarrow E, A \rightarrow B\}$$

Example 3: $F = \{ABC \rightarrow CD, BC \rightarrow D, A \rightarrow B, C \rightarrow D\}$

\Rightarrow Step 1: no common / ~~not~~ some variables in LHS

$$\begin{array}{l} A \rightarrow B \\ ABC \rightarrow CD \end{array} \quad \left. \begin{array}{l} A \rightarrow B \\ AC \rightarrow CD \end{array} \right\}$$

$$F_1 = \{A \rightarrow B, AC \rightarrow CD, BC \rightarrow D, C \rightarrow D\}$$

Step 3:

$$\left. \begin{array}{l} C \rightarrow D \\ AC \rightarrow CD \end{array} \right\} \Rightarrow C \rightarrow D$$

Phone: +91 844-844-0102

(Because C already determined D alone ∴ we don't need A separately. And

$$C \rightarrow CD \Rightarrow C \rightarrow C \text{ (festival)}$$

$$C \rightarrow D \text{ (expected)}$$

$$F_2 = \{ C \rightarrow D, A \rightarrow B, BC \rightarrow D \}$$

~~Again,~~ Again, $\left. \begin{array}{l} BC \rightarrow D \\ C \rightarrow D \end{array} \right\} \Rightarrow C \rightarrow D$

(C can determine D ∴ we don't need B hence we don't need rest of $C \rightarrow D$ as it is already there)

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

Hence, minimal cover of F is $F_3 = \{ A \rightarrow B, C \rightarrow D \}$

Example $F = \{ A \rightarrow BC, B \rightarrow AC, C \rightarrow AB \}$

⇒ We can write it as:

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

(using decomposition laws)

Since, we have $A \rightarrow B$ and $B \rightarrow C$

then using transitivity we can get $A \rightarrow C$

Hence, we can remove it.

Similarly, $B \rightarrow C$ and $C \rightarrow A$, then using transitivity again we can drop $B \rightarrow A$.

Mail: gatecse@appliedroots.com

Similarly, we can drop $C \rightarrow B$ using transitivity because we have $C \rightarrow A$ and $A \rightarrow B$.

So now FD set is $F_2 = \{A \rightarrow B, B \rightarrow C, C \rightarrow A\}$

(This FD set is cyclic FD's)

So F_2 is the minimal set/cover of F.

This problem is solved by minimizing the individual FD's only (Step 3) as we have not removed any FD.

Example 5: $F = \{A \rightarrow BC, CD \rightarrow E, E \rightarrow C, D \rightarrow EAH, ABH \rightarrow BD,$
 $DH \rightarrow BC\}$

\Rightarrow Step 1: No common LHS

$$\begin{array}{l} \downarrow \\ \{D \rightarrow E, \\ D \rightarrow A, \\ D \rightarrow H\} \end{array}$$

Step 2:

$$\begin{array}{l} \because D \rightarrow EAH \Rightarrow D \rightarrow E, D \rightarrow A, D \rightarrow H \\ \cancel{CD \rightarrow E} \\ \Downarrow \\ D \rightarrow EAH \end{array} \left. \begin{array}{l} \because D \text{ can} \\ \text{determine} \\ E \text{ uniquely} \\ \therefore \text{we don't} \\ \text{need } C \\ \text{and } H \text{ hence} \\ \text{we can drop} \\ \cancel{CD \rightarrow E} \end{array} \right\}$$

$$F_1 = \{A \rightarrow BC, E \rightarrow C, D \rightarrow EAH, ABH \rightarrow BD,$$

 $DH \rightarrow BC\}$

$$\therefore A \cancel{\rightarrow} H \rightarrow BD \Rightarrow AH \rightarrow D \quad (x_4 \rightarrow x_2 \Rightarrow y \rightarrow z)$$

$$F_2 = \{A \rightarrow BC, E \rightarrow C, D \rightarrow EAH, AH \rightarrow D, DH \rightarrow BC\}$$

$$\begin{array}{l} \therefore D \rightarrow EAH \Rightarrow D \rightarrow E, D \rightarrow A, D \rightarrow H \\ \cancel{DH \rightarrow BC} \end{array} \left. \begin{array}{l} DR^+ = \{D, H, E, A, C, B\} \\ \text{and } D \rightarrow H \\ \therefore \text{no need to have } D \end{array} \right\}$$

$$F_3 = \{ A \rightarrow BC, E \rightarrow C, D \rightarrow EAH, AH \rightarrow D \}$$

We cannot further reduce individual FD's further hence, $F^+ = F_3^+$. (F_3 is minimal cover of F)

16.6 Decomposition with Solved Problems:

Given a relation $R(A_1, A_2 \dots A_n)$ and set of FD's $F = \{\dots\}$, breaking R into further relations $R_1, R_2 \dots R_n$ such that all attributes of R get covered in decomposed relations. (Breaking a larger relation into smaller relation)

→ Properties of decomposition: (Properties to be held while performing decomposition)

→ No redundancy or minimal redundancy:

It means that we want to minimize the disk storage by avoiding storage of duplicate / redundant data. Hence, ~~we don't store data~~ decomposition should be done such that no duplicate / minimum duplicate data is stored.

→ Lossless Decomposition:

Decomposition should be done such that there is no loss of data while combining/joining them later. (lossless join)

Dependency preserving:

It means that while decomposing R into $R_1, R_2 \dots R_n$, all FD's in F should be preserved by $R_1, R_2 \dots R_n$.

(Normalization is a principled and systematic approach to decomposition)

To understand the concept of decomposition, let us consider an example below:

Cid	Pid	Cname	Pname	time
i ₁	j ₁	n ₁	p ₁	t ₁
i ₁	j ₂	n ₁	p ₂	t ₂
i ₂	j ₂	n ₂	p ₂	t ₃
i ₃	j ₃	n ₃	p ₃	t ₄

Fig:33 - Purchases

If we have only 1 relation like this, we will have lots of redundancy / repeated storage of lots of information.

Note:

There are several anomalies which we have to see in this concept of decomposition (anomalies which are present when we have only 1 relation)

* Update anomaly: If we want to change Pname from 'P1' to 'P2' in row 102, we need to update multiple cells which is time consuming.

* Insertion anomaly: If we want to insert data about a new customer having tuple values as follows. (Note that this is a new customer and has not purchased anything yet, therefore entry of this customer in the purchase table will be like)

(i4, NULL, n4, NULL, NULL)

Here, we are wasting a lot of space by having so many NULL values.

* deletion anomaly: Suppose we want to delete the information of product id Pid and Pname. If we delete the whole tuple, we will also lose customer information, so the only ~~update~~ alternate option is to update them to NULL.

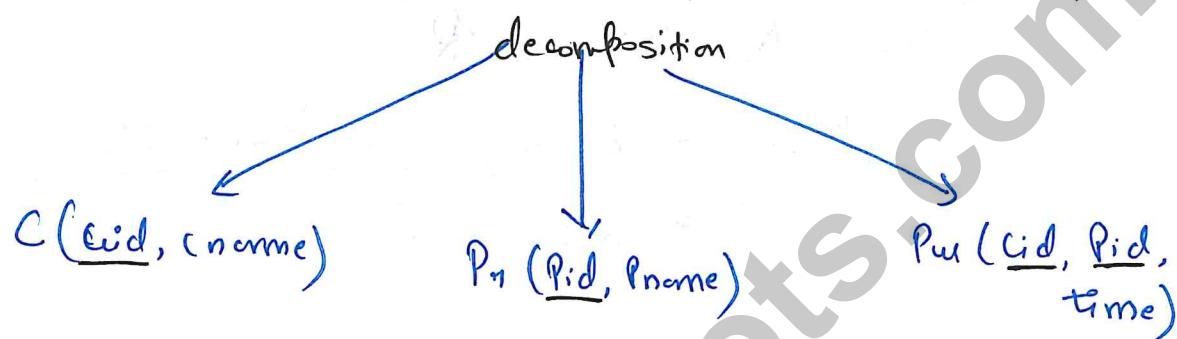
So deletion anomaly will result in ^{too} many NULL values in relation and updation of multiple cell to NULL values which is time consuming.

→ Eliminating Redundancy

We can eliminate redundancy by decomposition.

Examples

Purchases (cid, cname, pid, pname, time)



→ lossless Decomposition (lossless join decomposition):

If we decompose a relation R into R_1, R_2 and R_3 and $R_1 \bowtie R_2 \bowtie R_3 = R$ (Natural join) then the decomposition is called as lossless join decomposition, else it is lossy.

Example

R (A B C)

A	B	C
1	2	1
2	2	2
3	1	2

R

decomposed

$R_1 (A B)$

A	B
1	2
2	2
3	1

R_1

$R_2 (B C)$

B	C
2	1
2	2
1	2

R_2

here, $R_1 \bowtie R_2 =$

A	B	C
1	2	1
1	2	2
2	2	1
2	2	2
3	1	2

This was not
in original
table

Since, $R_1 \bowtie R_2 \supsetneq R$

hence, it is a lossy decomposition
($R_1 \bowtie R_2 \neq R$)

This didn't work because, here we are joining using the common attribute B and B is having duplicates, it means B that this common attribute is not a key or more specifically super key or it cannot determine a tuple uniquely neither in R_1 nor in R_2 . The decomposition would not be lossy if B is the key of in either of the relation.

So, if the common attribute is not a super key of any of the tables R_1 or R_2 than the decomposition is lossy.

Example 2: $R(A B C)$

Phone: +91 844-844-01952

From the instance of table R in example 1
the functional dependency which we could derive is
 $F = \{A \rightarrow B, A \rightarrow C\}$

\Rightarrow Given relation $R(A B C)$

A	B	C
1	2	1
2	2	2
3	1	2

R

decomposed

$R_1(A B)$

A	B
1	2
2	2
3	1

$R_2(A C)$

A	C
1	1
2	2
3	2

$R_1 \bowtie R_2$

A	B	C
1	2	1
2	2	2
3	1	2

$= R$

\Rightarrow This decomposition is lossless as
 $R_1 \bowtie R_2 = R$

In table R_1 , A is the superkey and A is also a superkey of R_2 .

Ques

So, if the common attribute in decomposition is a SR of either of the relation (resultant relations) then decom. is lossless.

Example 3: $R(A B C D)$

Phone: +91 844-844-0102



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

\rightarrow Given decomposition $R_1(A B C)$ and $R_2(C D)$ is lossless or not?

$\Rightarrow \because C \rightarrow D \therefore$ in relation $R_2(C D)$

C is the SK. Hence, the condition is satisfied and therefore, the decomposition is lossless.

\rightarrow Given decomposition $R_1(A B)$ and $R_2(C D)$ is lossless or not?

\Rightarrow Since, there is no common attribute in between R_1 and R_2 and therefore, we cannot join them. Hence, this decomposition is lossy.

\rightarrow Given decomposition $R_1(A B)$ and $R_2(B C)$ is lossless or not?

\Rightarrow It is also lossy because it doesn't even include the attribute D in any of the relation.

\rightarrow Given decomposition $R_1(A B)$, $R_2(B C)$ and $R_3(C D)$ is lossless or not?

$\Rightarrow \because B \rightarrow C \therefore B$ is SK to $R_2 \therefore$ we can combine R_1 and R_2

$\therefore C \rightarrow D \therefore C$ is SK to R_3 . Hence, we can combine R_2 with R_3 . Hence, it is lossless.

Mail: gatcse@appliedroots.com

NOTE

Algorithm to determine for lossy / lossless decomposition

Phone: +91 844-814-0102

① R, F

$\hookrightarrow R_1, R_2 \dots R_n$ is lossless if

$$\textcircled{a} \quad R \cup R_2 \cup \dots \cup R_n = R$$

\textcircled{b} Merging R_i and R_j into R_{ij}

$$ij \quad \textcircled{a} \quad R_i \cap R_j \neq \emptyset$$

$$\textcircled{b} \quad R_i \cap R_j \rightarrow R_i, \text{ or}$$

$$R_i \cap R_j \rightarrow R_j$$

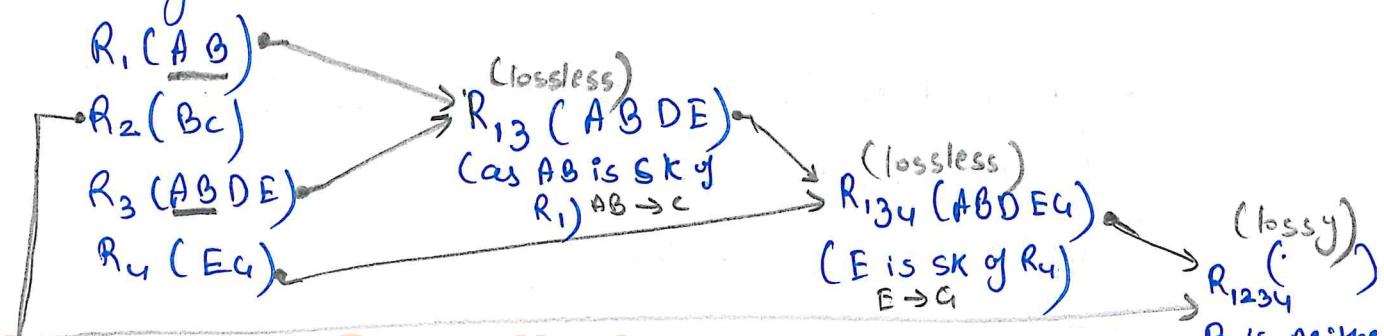
repeat
till we get
original R

Example: $R(A B C D E G)$

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

Given decomposition $R_1(AB), R_2(BC), R_3(ABDE)$
 $R_4(EG)$ is lossless or not?

Given the relations R_1, R_2, R_3 and R_4 union of all attributes are present in R . So, \textcircled{a} condition is satisfied.



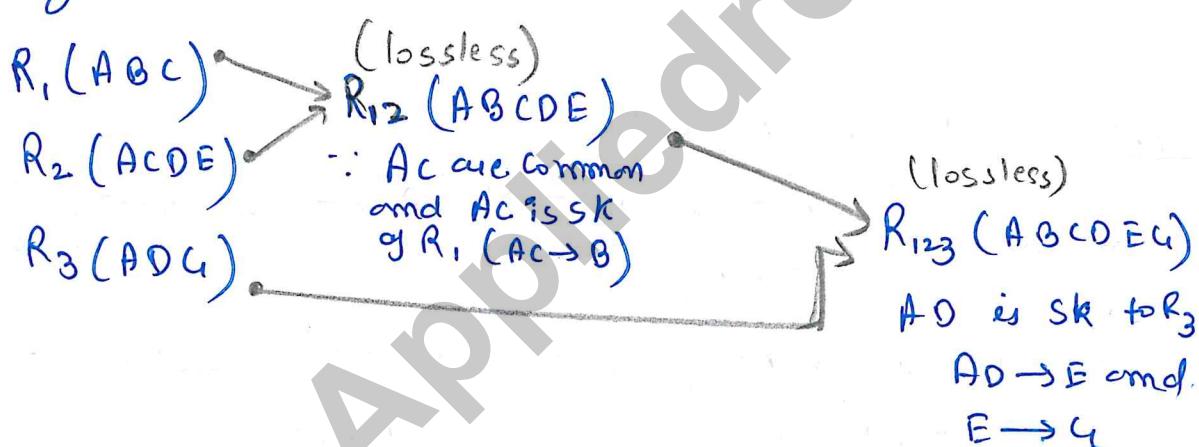
Hence, this whole decomposition is lossy.

We can try multiple ways to merge them but they will all return us a lossy decomposition.

→ Given decomposition $R_1(ABC)$, $R_2(ACDE)$, $R_3(ADG)$ is lossless or not?

⇒ Since, decomposition covered all attributes \therefore first condition is satisfied.

Now to check second condition we need to merge them

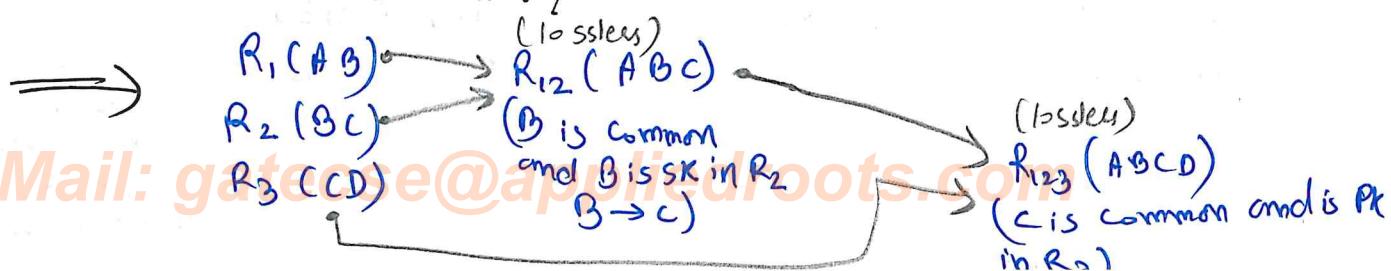


Hence, this decomposition is lossless

Example: $R(ABCD)$

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

→ Given decomposition $R_1(AB)$, $R_2(BC)$, $R_3(CD)$ is lossless or not?

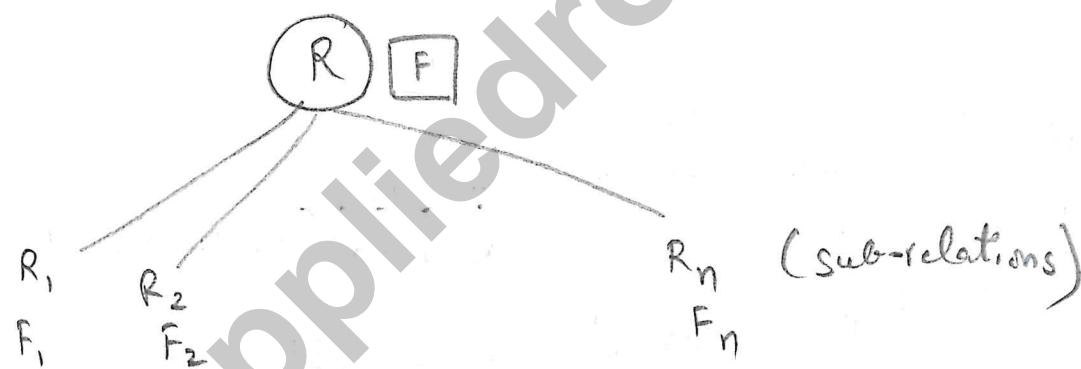


Hence, this decomposition is lossless.

16.7 Dependency preserving decomposition + Solved examples:

→ Dependency preserving decomposition:

When we try to decompose the relation R into $R_1, R_2 \dots R_n$, we also try to decompose the FD set F, into other smaller set of FD's $F_1, F_2 \dots F_n$.



So, decomposition is said to be dependency preserving if $F_1 \cup F_2 \dots \cup F_n = F'$ and $F' = F$ i.e two FD set F' and F should be equivalent. \leftarrow

$$\text{i.e. } (F_1 \cup F_2 \dots \cup F_n)^+ = F^+$$

(Recover F only using F_i 's)

Example f: $R(A B C D)$

Phone: +91 844-844-0102

APPLIED
ROOTS

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

Given the decomposition of R as $R_1(A B)$, $R_2(B C)$

and $R_3(C D)$ is dependency preserving or not?

\Rightarrow Step 1: Identify all non-trivial FD's using F for each sub relations.

$$R_1(A B) \rightarrow F_1 = \{ A \rightarrow B, B \rightarrow A \}$$

$$R_2(B C) \rightarrow F_2 = \{ B \rightarrow C, C \rightarrow B \}$$

$$R_3(C D) \rightarrow F_3 = \{ C \rightarrow D, D \rightarrow C \}$$

$$\left. \begin{array}{l} \\ \\ \end{array} \right\} F'$$

$$A \rightarrow B$$

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

Hence, we can add $B \rightarrow A$ (additional FD)

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

it is additional FD in R_2

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

it is additional FD in R_3

Step 2: Determine if $F_1 \cup F_2 \cup F_3 = F$ or not?

Here, we need to check whether we can get all FD's in F' using F or not or viceversa

$$F' \text{ covers} = \{ A \rightarrow B, B \rightarrow C, C \rightarrow D \}$$

Mail: gatcse@appliedroots.com
 $D \rightarrow A$ is not yet covered but we can infer it using transitivity. Hence, we can recover F using F' . It is DP.

Example 2: $R(AB \cup DEF)$

$R(AB \subset DEF)$

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

Given decomposition $R_1(A B)$, $R_2(C D E)$, $R_3(E F)$ is

DP (dependency preserving) or not?

\Rightarrow Step 3:

$$R_1(A B) \rightarrow f_i = \{ A \xrightarrow{\quad} B, B \xrightarrow{\quad} A \} \quad A^+ = A \text{ and } B^+ = B \\ \text{Hence, no FD is valid}$$

$A^t = A$ and $B^t = B$
Hence, no FD is valid

$$R_2(CDE) \rightarrow F_2 = \{ \underset{C \rightarrow D}{\text{C} \rightarrow D}, C \rightarrow E, D \rightarrow E, \} \quad C^+ = \{ C, D, E, F \} \quad O^+ = \{ D, E, F \}$$

$$R_3(EF) \rightarrow F_3 = \{ \overset{CD \rightarrow E}{E \rightarrow F}, \overset{E^+ = \{E, F\}}{\cancel{F \rightarrow E}} \} \quad , \quad \begin{cases} CD^+ = \{C, D, E, F\} \\ E^+ = \{E, F\} \end{cases}, F^+ = F$$

Step 2: Check $F_1 \cup F_2 \cup F_3 = F$

On checking F and F' we can "never get" $AB \rightarrow c$ and $AB \rightarrow D$
 \therefore This decomposition is not dependency preserving.

Example 3: $R(A B C D E)$

$$F = \{A \rightarrow BC, CD \rightarrow E, B \rightarrow D, E \rightarrow A\}$$

Given decomposition $R_1(ABC)$ and $R_2(AOE)$ is DP or not?

1

Steps:

$$R_1(A \cup B \cup C) \Rightarrow F_1 = \{A \rightarrow BC, A \rightarrow B, A \rightarrow C, AB \rightarrow C\}$$

$\hookrightarrow A, BC \rightarrow B, AC \rightarrow B$

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

$$R_2(ADE) \rightarrow F_2 = \{ E \rightarrow A, A \rightarrow E, A \rightarrow D, E \rightarrow D \} \quad A^+ = \{A, C, D, E\}$$

$AE \rightarrow D, AD \rightarrow E, DE \rightarrow A$

Step 2: $\therefore A \rightarrow B C$ is covered in R.

and $\text{CD} \rightarrow \text{E}$ $\text{E} \rightarrow \text{A}$ is covered in R_2

Belt for $B \rightarrow D$ and $CD \rightarrow E$ we need to check.

APPLIED ROOTS

$$B_{F_1, F_2}^+ = \{B\}$$

(B^+ under
 F_1 and F_2)

$$C_{D_{F_1, F_2}}^+ = \{C, D\}$$

$\therefore B \rightarrow D$ and $CD \rightarrow E$ is not preserved
Hence, this decomposition is DP.

Example 4: $R(ABCDEFHI)$

$$F = \{AB \rightarrow C, A \rightarrow DE, B \rightarrow F, F \rightarrow GH, D \rightarrow IJ\}$$

Decomposition is given as $R_1(ABC)$, $R_2(ADE)$,
 $R_3(BF)$, $R_4(FGH)$, $R_5(DIJ)$ is DP or not?



Step 1:

$$R_1(ABC) \rightarrow F_1 = \{AB \rightarrow C, \dots\}$$

$$R_2(ADE) \rightarrow F_2 = \{A \rightarrow DE, \dots\}$$

$$R_3(BF) \rightarrow F_3 = \{B \rightarrow F, \dots\}$$

$$R_4(FGH) \rightarrow F_4 = \{F \rightarrow GH, \dots\}$$

$$R_5(DIJ) \rightarrow F_5 = \{D \rightarrow IJ, \dots\}$$

here, each of the ~~relations~~ ~~are~~ decomposition are done such that for every FD there is a separate relation. Since, there is no need to look for the additional non-trivial dependency as they are completely preserved. Hence, it is DP. It is by design DP.

Suppose the decomposition is like: **Phone: +91 844-844-0102**

$$R_1(A B C D E) \rightarrow F_1 = \{AB \rightarrow C, A \rightarrow DE \dots\}$$

$$R_2(B F G H) \rightarrow F_2 = \{B \rightarrow F, F \rightarrow GH \dots\}$$

$$R_3(D I J) \rightarrow F_3 = \{D \rightarrow IJ \dots\}$$

It is also DP.

NORMALIZATION

17.1 Introduction to Normalization + INF and 2NF + examples :

- * Principled approach to perform decomposition in order to ensure
 - No or minimum redundancy
 - lossless join
 - Dependency preserving
 - Reduce anomalies
- * Covers normal forms like 1NF, 2NF, 3NF, BCNF, 4NF etc.

Edgar
Codd
came up
with
Principled
approach
Normalization

Note that, every relation in 1NF is in 2NF, every relation in 2NF is also in 3NF and so on. So, every higher normal form implies lower normal form.

\rightarrow $\in 1NF$:

Phone: +91 844-844-0102

Each attribute in a relation is atomic (single valued or not multivalued)

Example 1: Relation R:

Cid	Purchased
i ₁	P ₁ , P ₂
i ₂	P ₂ , P ₃

This attribute is not atomic and $\therefore R$ is not in 1NF

Relation R':

Cid	Purchased
i ₁	P ₁
	P ₂
i ₂	P ₂
	P ₃

Redundancy

decomposing it, so that it could be in 1NF

it is in 1NF
But here PK is {Cid, Purchased}

One of the problems with 1NF is that it introduce redundancy. But is the foundation of normal forms as we cannot have multivalued attributes or multiple values in a cell.

Example 2: R(A B C D) where D is multivalued attribute. FD, F = {A \rightarrow B, B \rightarrow C} what are possible relations in which are in 1NF?

\Rightarrow Possible CK are:

\because D is not in FD's \therefore D should be in every CK.

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

Hence, $CK = \{AD\}$. Since, AD is a Candidate Key
~~so~~ therefore D could not be multivalued. The relation will be

$R^*(A, D, BC)$ where A D is Key D will be broken into

$R_1(A, D)$, $R_2(A, BC)$

Here D

will be

splitted

like shown in example 1

→ Partial dependency (PD) and Redundancy =

Partial dependency is given as ⇒

FD : Proper subset of a candidate key → non prime-attribute

Example :

R

Cid	Cname	Pid
C1	n1	P1
C1	n2	P2
C2	n2	P2
C2	n2	P3

Redundancy due to partial dependency

Subset of CK ↑

FD, F = {Cid → Cname}

non-prime attr below

here $CK = \{(cid, pid)\}$

Prime-attribute = {cid, pid}

non prime-attribute = {cname}

The FD given in set F is a Partial dependency.

\rightarrow 2NF:

Phone: +91 844-844-0102



A relation R is in 2NF iff

- ① R is in 1NF (all attributes are atomic)
- ② R has no partial dependency.

2NF avoids the redundancy caused due to partial dependency.

Example: R (Cid, Pid, cname)

$$FD = \{ Cid \rightarrow cname, Cid, Pid \rightarrow cname \}$$

Partial dependency

R is already in 1NF because all attributes are atomic). Now if we decompose

R (Cid, Pid, cname)

R₁ (Cid, Pid)

$$F_1 = \{ Cid, Pid \rightarrow \}$$

(R₁ is in 2NF)

$$CK = \{ (Cid, Pid) \}$$

Hence, PD

R₂ (Cid, cname)

$$F_2 = \{ Cid \rightarrow cname \}$$

(R₂ is in 2NF)

$$CK = \{ Cid \}$$

Hence, no PD

Hence, decomposition of R into R₁ and R₂ is in 2NF and is a lossless ; dependency preserving decomposition .

Example:

$R(ABCDEF)$ is in 1NF

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

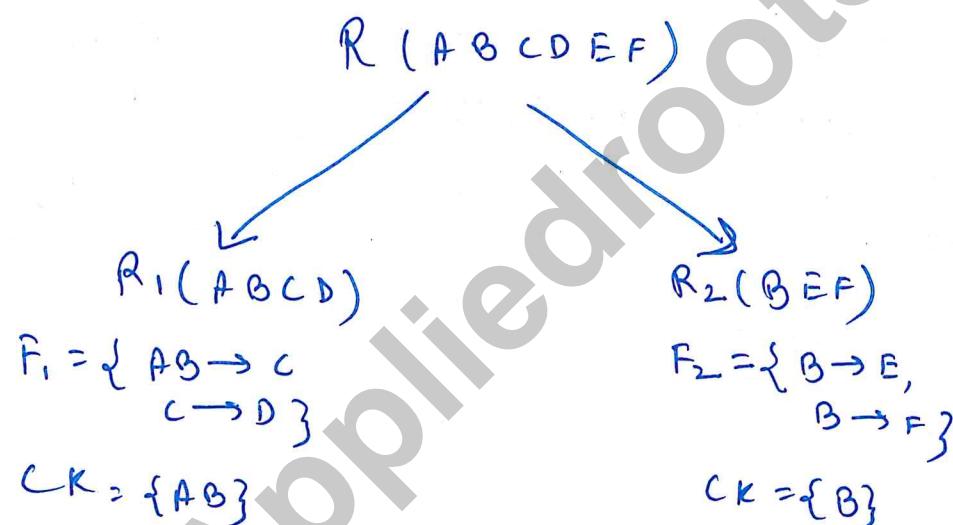
Is it is 2NF. If not then decompose.

$$\Rightarrow CK = \{AB\} \quad \because AB^+ = A, B, C, D, E, F$$

Prime attribute = {A, B}

here $B \rightarrow E$ and $B \rightarrow F$ are Partial dependencies

Therefore we need to decompose it :



\therefore All dependency's are preserved and
it is a lossless decomposition as common attribute
 B is SK of R_2 . Both R , and R_2 are in 2NF
as there are no PD.

Example:

$R(ABCDEFG)$ is in 1 NF

$$FD = \{AB \rightarrow C, B \rightarrow F, A \rightarrow E, E \rightarrow G\}$$

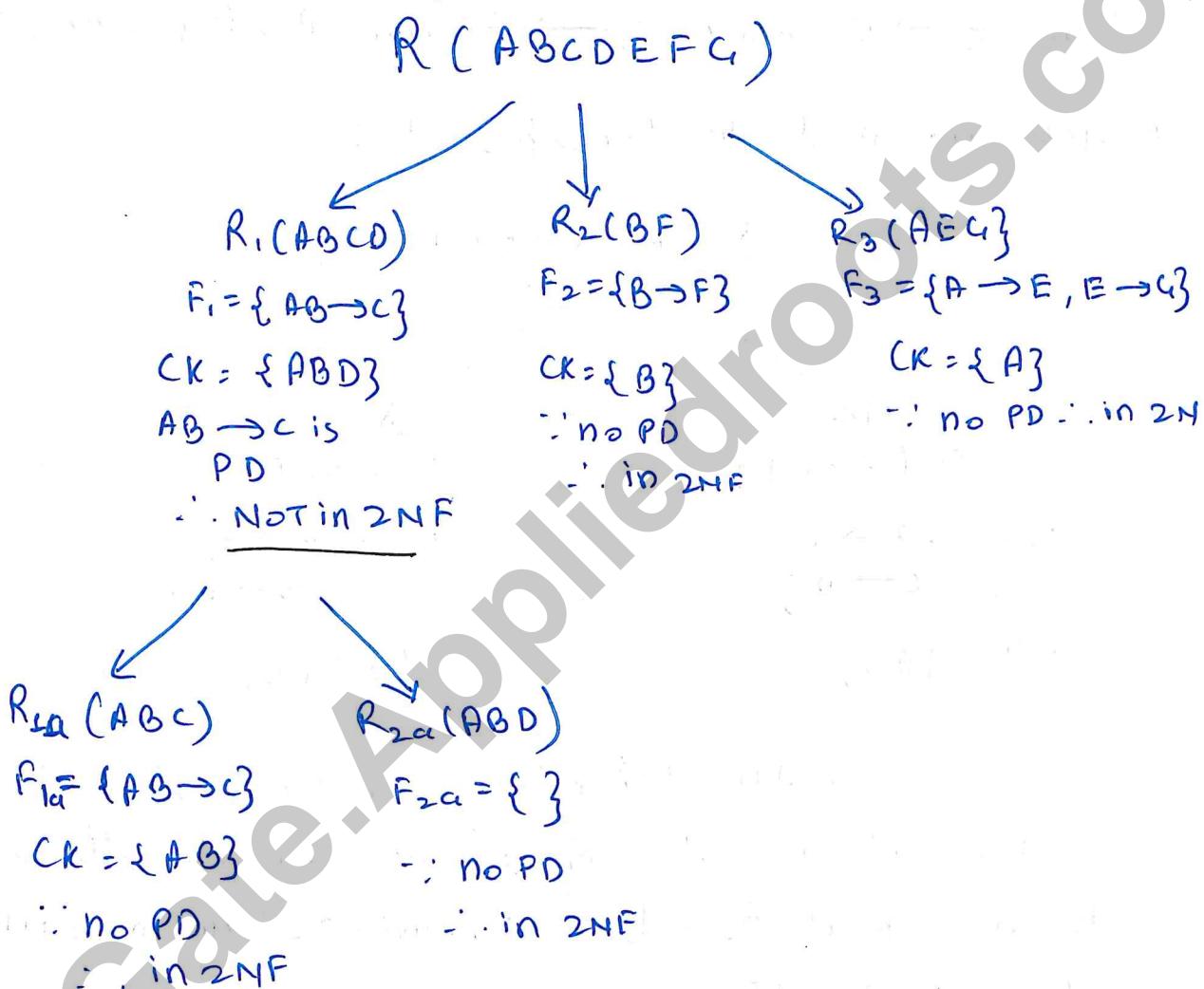
is R in 2NF? If not then decompose it.

here, $CK = \{ ABD \}$

Prime attribute = {A, B, D}

Here, $A \rightarrow C$, $B \rightarrow F$, $A \rightarrow E$ are PD $\therefore R$ is not in 2NF

We need to decompose R:



NOTE - 2NF does not remove all redundancies.
 It eliminates only redundancy due to partial dependencies.

17.2 3NF and BCNF:

Tournament winners:

<u>Tournament</u>	<u>Year</u>	Winner	Winner date of birth
Indiana Invitational	1998	Al Fred	21 July 1975
Cleveland open	1999	Bob Albert	28 Sept 1968
Des Moines	1999	Al Fred	21 July 1975
Indiana Invitational	1999	Chip Masterson	14 March 1977

Tournament winners:

<u>Tournament</u>	<u>Year</u>	Winner
Indiana Invitational	1998	Al Fred
Cleveland open	1999	Bob Albert
Des Moines	1999	Al Fred
Indiana Invitational	1999	Chip Masterson

Winner Date of birth:

<u>Winner</u>	<u>Date of Birth</u>
Chip Masterson	14 March 1977
Al Fred	21 July 1975
Bob Albert	28 Sept 1968

→ Transitive Dependencies

Considering Tournament as T Phone: +91 844-844-0102

Year as Y

Winner as W

Winner date of birth as WD

Given, FD set $F = \{ TY \rightarrow WD, W \rightarrow D \}$

In the relations given, there is some redundancy ^{in R}, though it is already in 2NF.

$CK = \{ TY \}$, and there is no PD. Hence, it is an 2NF.

Since, prime attributes are $= \{ T, Y \}$

Non-prime attributes are $= \{ W, D \}$

here $W \rightarrow D$ is like non-prime attribute determining another non-prime attribute. The redundancy in R is due to this dependency. Such dependencies are called Transitive dependency.

→ 3NF's

A relation R is in 3NF iff

(a) R is in 2NF, and

(b) All non-prime attribute must depend on a key (sk)

Better way to understand

{ 'OR' for all non-trivial FD's $X \rightarrow Y$
either X should be a sk or Y is a prime attribute.

'OR' No non-prime attribute is transitively determined independent on any key

'OR' No non-prime attribute depends on other non-prime attribute. such dependences are called Transitive dependency

Example:

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

Phone: +91 844844102
SK prime attribute

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

Is R in 3NF? If not decompose it.

Given the FD set F, the CK is

$$CK = \{ AB \}$$

$$\text{prime-attribute} = \{ A, B \}$$

here, $B \rightarrow F$ is a partial dependency
and $C \rightarrow D, D \rightarrow E$ is a transitive dependency

Therefore, R is neither in 2NF nor in 3NF.

If the decomposition is

By design it is DP and lossless

$$\left\{ \begin{array}{l} R_1(ABC) \rightarrow F_1 = \{ AB \rightarrow C \} \\ R_2(CD) \rightarrow F_2 = \{ C \rightarrow D \} \\ R_3(DE) \rightarrow F_3 = \{ D \rightarrow E \} \\ R_4(BF) \rightarrow F_4 = \{ B \rightarrow F \} \end{array} \right.$$

Now they are in 2NF as well as in 3NF

Hence, R_1, R_2, R_3, R_4 are in 3NF. No other decomposition is possible to have it in 3NF.

→ ↴ ← BCNF :

It stands for Boyce-Codd Normal form. It is a stricter form of 3NF. Given a relation R and FD set F, every FD $X \rightarrow Y$ is in BCNF such that

(a) $X \supseteq Y$ (trivial)

(b) X is a superkey.

Every BCNF is in 3NF. Note that there are very few real world examples of 3NF and not in BCNF.

~~NOTE~~ - Some computer scientist Bernstein said in their research paper.

Bernstein in 1979 said in their research paper.

Given the following type of ~~partial~~ relation and F)

$R(ABC)$

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

We cannot obtain BCNF which is lossless, dependency preserving and ~~BCNF~~ in BCNF. At any one condition will be compromised.

Also it is noted that, having a lossless join is a must / necessary condition while this is not the case with DP.

We try to preserve all functional dependency's but it is not a must condition.

Proof

$R(ABC)$, $F = \{ AB \rightarrow C, C \rightarrow B \}$

$$\begin{aligned} R_1(A\bar{C}) &= \{ \} \\ R_2(BC) &= \{ C \rightarrow B \} \end{aligned} \quad \left. \begin{array}{l} \text{Both are in} \\ \text{BCNF} \end{array} \right\}$$

R_1 and R_2 are in BCNF and will have a lossless join as C is sk of R_2 . But it could not preserve $AB \rightarrow C$. Hence, we need to compromise with the dependency preserving condition and then we could not preserve $AB \rightarrow C$.

Example: $R(A B C D)$

$$F = \{ AB \rightarrow CD, D \rightarrow A \}$$

Is it in BCNF? If not then decompose it



$$CK = \{ AB, BD \}$$

$$\text{Prime attribute} = \{ A, B, D \}$$

For the given FD set F, R is in 3NF but not in BCNF because given,

decomposition1: \times

$$R_1(ABC) \Rightarrow F_1 = \{ \}$$

$$R_2(AD) \rightarrow F_2 = \{ D \rightarrow A \}$$

} Both are in BCNF
But they are neither
DP nor having
lossless join

decomposition2: \therefore

$$R_1(BCD) \rightarrow F_1 = \{ \}$$

$$R_2(AD) \rightarrow F_2 = \{ D \rightarrow A \}$$

} Both are in BCNF
and have lossless join
but it is not
DP as it could
not preserve
 $AB \rightarrow CD$

17.3 Solved Problems + Properties of Normal Forms

Example 1: $R(A B C D E F G H I J)$

$$F = \{ A B \rightarrow C, A \rightarrow D E, B \rightarrow F, F \rightarrow G H, D \rightarrow I J \}$$

Check whether the given relation and FD set is in 2NF, 3NF, BCNF.

\Rightarrow Candidate Keys = { AB }

Prime attributes = { A, B }

* * Check for 2NF * * (No PD is part of key \rightarrow non-prime attribute)

Here, $A \rightarrow D E$ and $B \rightarrow F$ are Partial dependencies

∴ We will decompose R such that, taking the dependency which are violating 2NF and then break/make a separate relation for them.

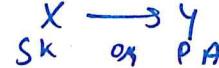
$R_1 (A D E I J) \rightarrow F_1 = \{ A \rightarrow D E, D \rightarrow I J \} // A \text{ is key}$

$R_2 (B F G H) \rightarrow F_2 = \{ B \rightarrow F, F \rightarrow G H \} // B \text{ is key}$

$R_3 (A B C) \rightarrow F_3 = \{ A B \rightarrow C \} // A B \text{ is key}$

Since, there is not PD left ∴ R_1, R_2, R_3 are in 2NF now.

* * Check for 3NF * * **(Phone 191 844-024-0102)**



21

Among R_1, R_2, R_3 , R_2 has the relation R_1 and R_2 are having transitive dependency
 $D \rightarrow IJ$ and $F \rightarrow GH$ are TD

The decomposition will be:

By design in BCNF

$R_1(ADE) \rightarrow F_1 = \{A \rightarrow DE\}$	3NF / BCNF
$R_2(BF) \rightarrow F_2 = \{B \rightarrow F\}$	3NF / BCNF
$R_3(FGH) \rightarrow F_3 = \{F \rightarrow GH\}$	3NF / BCNF
$R_4(DIJ) \rightarrow F_4 = \{D \rightarrow IJ\}$	3NF / BCNF
$R_5(ABC) \rightarrow F_5 = \{AB \rightarrow C\}$	3NF / BCNF

While decomposing it into 3NF, we also got it into BCNF
Hence, R_1, R_2, R_3, R_4, R_5 are DP, lossless and BCNF decomposition

Example 2: $R(ABCD EFGH)$

$$F = \{AB \rightarrow C, AC \rightarrow B, AD \rightarrow E, B \rightarrow D, BC \rightarrow A, E \rightarrow G\}$$

Check for 2NF, 3NF and BCNF.

$$\Rightarrow CK = \{FHAB, FHAC, FHBC\}$$

Prime attribute = {A, B, C, F, H}

Ques

NOTE: In practice people directly convert relations into BCNF rather than converting them into 2NF and 3NF first.

Because everything which is in BCNF will automatically be in 2NF and 3NF.

* * Check for 2NF **

Phone: +91 844-844-0102



In FD set F, $B \rightarrow D$ is the functional dependency
∴ Decomposition will be:

$$R_1(BD) \rightarrow F_1 = \{B \rightarrow D\}$$

$$\times R_2(ABCEFH) \rightarrow F_2 = \{AB \rightarrow C, AC \rightarrow B, BC \rightarrow A, E \rightarrow G\}$$

But we cannot
Put $AD \rightarrow E$
∴ we cannot have this R_2

$$R_2(ABC FH) \rightarrow F_2 = \{AB \rightarrow C, AC \rightarrow B, BC \rightarrow A\}$$

$$R_3(ADEG) \rightarrow F_3 = \{AD \rightarrow E, E \rightarrow G\}$$

∴ Now R_1, R_2, R_3 is in 2NF

* Check for 3NF and BCNF **

In the above decomposition R_3 violates 3NF

∴ R_1 and R_2 will be same

$$R_3(ADE) \rightarrow F_3 = \{AD \rightarrow E\} \quad 3NF/BCNF$$

$$R_4(EH) \rightarrow F_4 = \{E \rightarrow H\} \quad 3NF/BCNF$$

Now R_1, R_2, R_3, R_4 are in 3NF

But R_2 is not in BCNF because $CK = FHAB, FHBC, FHAC$

Special case { ∴ $R_{2a}(ABC) \rightarrow F_{2a} = \{AB \rightarrow C, AC \rightarrow B, BC \rightarrow A\}$

$$R_{2b}(ABFH \setminus ACFH \setminus BCFH) \rightarrow F_{2b} = \{\}$$

Mail: gatecse@appliedroots.com

Now, $R_1, R_{2a}, R_{2b}, R_3, R_4$ are in BCNF.

Example 3: $R(ABCDE)$

$$F = \{ A \rightarrow BC, CD \rightarrow E, B \rightarrow D, E \rightarrow A \}$$

Check R for 3NF and BCNF.

$$\Rightarrow CK = \{ A, E, CD, BC \}$$

Prime attributes = { A, B, C, D, E }

$$R_1 (BD) \rightarrow F_1 = \{ B \rightarrow D \}$$

$$R_2 (AEC) \rightarrow F_2 = \{ A \rightarrow BC \}$$

$$R_3 (EA) \rightarrow F_3 = \{ E \rightarrow A \}$$

3NF / BCNF

But $CD \rightarrow E$ is not preserved

Hence R_1, R_2, R_3 are in BCNF and have lossless join but is not dependency preserving.

→ Properties of Normal forms:

- * If every CK is made of a single attribute (Simple Key) then Relation R and FD set F is in 2NF

Ex- $R(ABCDE)$

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

$$CK = \{ A, B, C \}$$

// Since all of them are simple keys therefore, it is by default in 2NF

Because, we can check there is no PD. (Part of CK → non prime)

Mail: gatcse@appliedroots.com

* If every attribute is a prime attribute then relation R, and FD set F is in 3NF

Phone: +91 844-844-0102

* If relation R is in 3NF and all candidate keys are simple than relation R and FD set P is in BCNF.

(if every ck is simple than every prime attribute are ck and R is in 3NF which means $x \rightarrow y$ is satisfied \hookrightarrow_{SK})

(if x is a SK than BCNF and if x is not SK then will not arrive)

* A relation with only 2 attributes (Binary relation) is always in BCNF.

Ex - Possible cases for R(AB) are:

- ① $\{ \}$; BCNF
- ② $\{ A \rightarrow B \}$; BCNF
- ③ $\{ B \rightarrow A \}$; BCNF
- ④ $\{ A \rightarrow B, B \rightarrow A \}$; BCNF
- ⑤ $\{ AB \rightarrow AB \}$; BCNF

17.4 Multi-valued dependencies and 4NF:

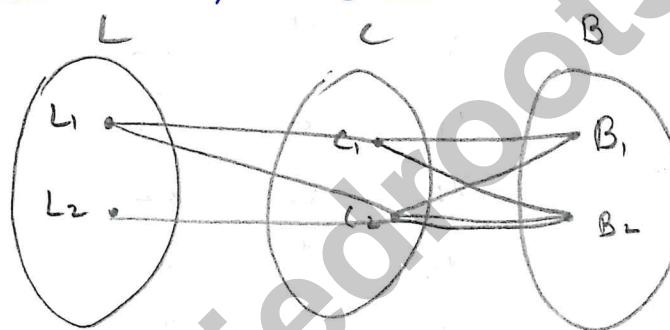
Phone: +91 844-844-0102

R:

Lecture (L)	Course (C)	Book (B)
L_1	C_1 / C_2	B_1 / B_2
L_2	C_2	B_2

Fig: 35 Wikipedia example

As we can see in fig 35 Course (C) and Book (B) are multivalued attribute. Hence, it is not in 4NF.



It could be written as

R'	L	C	B
	L_1	C_1	B_1
	L_1	C_1	B_2
	L_1	C_2	B_1
	L_1	C_2	B_2
	L_2		

4NF of R

4NF address this type of redundancy.

So, now R' is now in 4NF. But here no single attribute could be the key.

So, CK = {LCB}

There are no non-trivial FD in R' . ie $F = \emptyset$.

Therefore, R' is already in BCNF. But there is lot of redundancy in R' even though it is in BCNF.

Note that, this redundancy occurs because of multivalued attribute and ~~it~~ this type of redundancy is addressed by 4NF.

This type of redundancy could not be addressed by functional dependencies, therefore, the concept of Multivalued dependency come into the picture (MVD). They are called so because the origin of the problem is multivalued attributes.

→ Multivalued Dependency (MVD):

X	Y	Z
x_1	y_1	z_1
x_1	y_1	z_2
x_2	y_2	z_1
x_2	y_2	z_2

MVD is represented either using ' $\rightarrow\rightarrow$ ' or ' $\rightarrow\rightarrow\rightarrow$ '.

For a relation $R(XYZ)$, $Z=R-(X\cup Y)$ then
 X multi-determines Y ($X \rightarrow\rightarrow Y$).

X multi-determines Y if we have $Z=R-(X\cup Y)$,
and

$$\textcircled{1} \quad t_1 \cdot x = t_2 \cdot x = t_3 \cdot x = t_4 \cdot x$$

$$\textcircled{2} \quad t_1 \cdot y = t_2 \cdot y \text{ and } t_3 \cdot y = t_4 \cdot y$$

$$\textcircled{3} \quad t_1 \cdot z = t_3 \cdot z \text{ and } t_2 \cdot z = t_4 \cdot z$$

Properties of MVDs

- * If $x \rightarrow y$ then $x \rightarrow z$ when $z = R - (x \cup y)$
(Complement property)
- * Trivial MVD:
 - (a) if $X \supseteq Y$ then $X \rightarrow Y$
 - (b) if $X \cup Y = R$ then $X \rightarrow Y$
 * We need atleast 3 attribute for MVD to exist.
 * We need atleast 2 attribute for non-trivial FD to exist
- * If $x \rightarrow y$ and $Z \supseteq W$ then,
 $xZ \rightarrow YW$ (Augmentation)
- * If $x \rightarrow y$ and $y \rightarrow z$ then $x \rightarrow (z-y)$
(Transitivity)
- * If ~~not~~ $x \rightarrow y$ then $x \rightarrow y$ (every FD is MVD
but vice versa is not true)
(Replication)
- * MVD do not obey:
 - ⇒ if $x \rightarrow yz$ then $x \rightarrow y$ and $x \rightarrow z$ X
(it is not obeyed in MVD but it is obeyed in FD)
 - ⇒ if $x \rightarrow y$ and $x \rightarrow z \Rightarrow x \rightarrow yz$ X
(FD obey this ie $x \rightarrow y$ and $x \rightarrow z \Rightarrow x \rightarrow yz$)

→ 4NF:

Phone: +91 844-844-0102

Q7 is extension to BCNF. A relation R is in 4NF iff

(1) R, F is in BCNF (SK's are determined using set of FD's F)

(2) All non-trivial MVD's of $x \rightarrow y$. X should be SK and it should be determined using FD and not using MVD's

On considering relation given in fig: 36, ie

R	L	C	B	R_1	R_2
	L_1	C_1	B_1		
	L_1	C_1	B_2		
	L_1	C_2	B_1		
	L_1	C_2	B_2		
	L_2	C_2	B_2		

here, $F = \{ \}$

$$MVD = \{ L \rightarrow C, L \rightarrow B \}$$

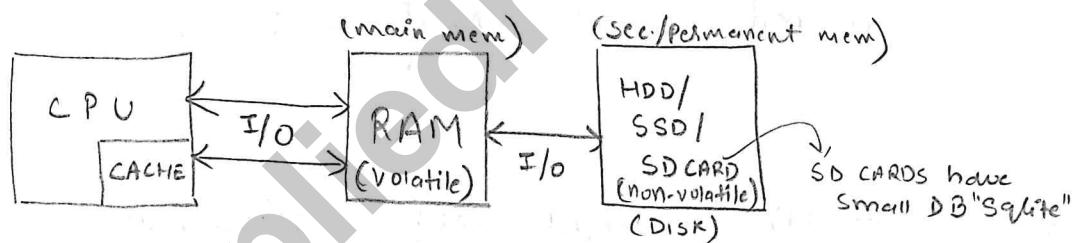
We can decompose R into 2 R_1 and R_2 and both R_1 and R_2 do not have any non-trivial MVD's

TRANSACTION AND

CONCURRENCY CONTROL

18.1 Model of a computer for transaction and Concurrency:

→ model of computer:



SQL query gets executed in CPU on database which was stored in Disk.

Failures possible while executing an SQL query:

- Power failure
- Disk failure
- Process failure

→ Transaction :

Transaction is a logical set of operations to perform a task.

Example: let us consider a bank transaction example:

- ① Read balance for A (\rightarrow Disk I/O operation)

- Phone: +91 844-044-0702
- ② If $A.\text{balance} \geq 500$ then $A.\text{balance} = A.\text{balance} + 500$
~~write A~~ (It will be performed in CPU and updated in RAM)
- ③ Write $A.\text{balance}$ ($\xrightarrow{\text{RAM or Disk}} \text{I/O}$ operation)
- ④ Read $B.\text{balance}$ ($\rightarrow \text{Disk I/O}$)
- ⑤ $B.\text{balance} = B.\text{balance} + 500$ ($\rightarrow \text{RAM operation}$)
- ⑥ Write $B.\text{balance}$ ($\xrightarrow{\text{RAM or Disk}} \text{I/O}$)
- ⑦ Commit (Commit stores data to non-volatile storage)

Each of these are logical operation.

Note - read data A \Rightarrow corresponds to R(A)
 Write data A \Rightarrow corresponds to W(A)
 Commit \Rightarrow corresponds to C

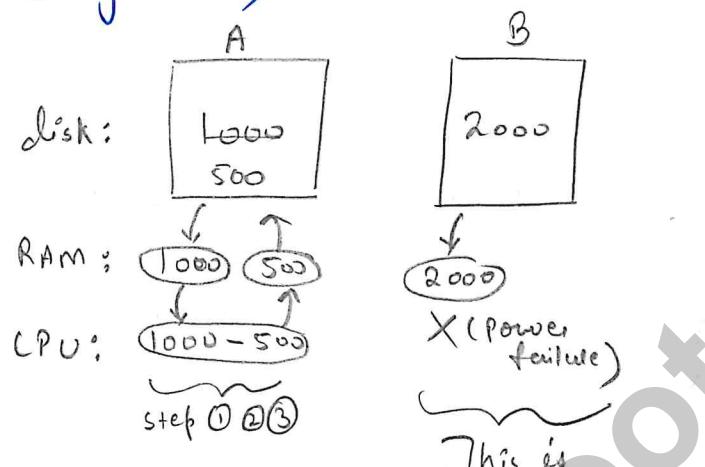
18.2 ACID Properties and Concurrency:

\rightarrow ACID Properties:

- A \rightarrow Atomicity
- C \rightarrow Consistency
- I \rightarrow Isolation
- D \rightarrow Durability.

→ - Atomicity:

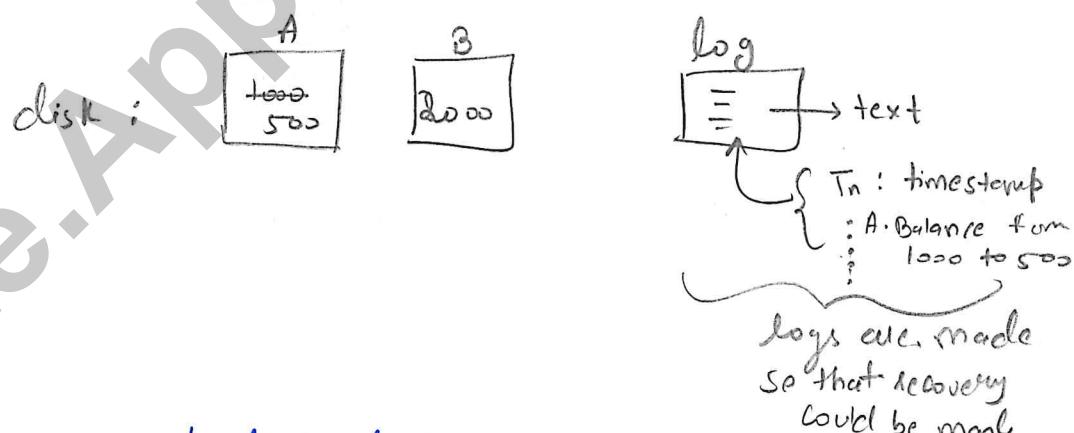
A transaction must execute all of its operations and commit successfully OR None of the operations should be executed (Should not have partial completion of operation)



This is partial completion
So a transaction should be atomic

If there is partial execution, then recovery should be done. One of the recovery technique used is log-based

recovery



So, during any kind of failure

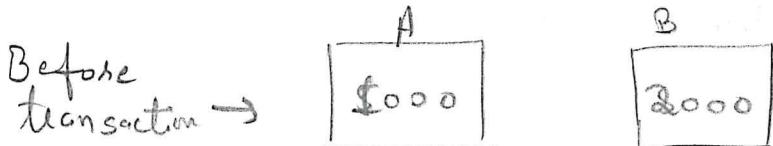
- read log file
- Observe partial completion
- Revert the transaction

So, Atomicity is a must property which tells a transaction should be either completely done or not done.

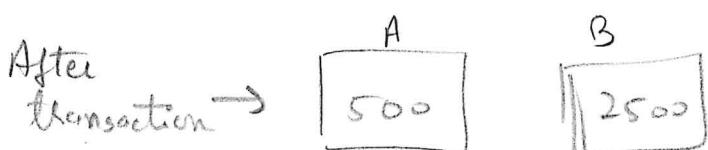
Consistency: The database state / state of system must be always consistent.

Phone: +91 844-844-0102

Ex- Transfer 500 from Acc. A to B (inter bank transfer)



$$A + B = 3000 \text{ (before transaction)}$$



$$A + B = 3000 \text{ (After transaction)}$$

If the system follows the following then that it will always remain in consistent state

Consistent $\xrightarrow{\text{follows}}$ Atomicity, Isolation $\xrightarrow{\text{and}}$ Consistent
Durability

We can enforce consistency using CHECK in sql

NOTE-

Concurrency:

In real world DB, tons of query are asked at same time (Concurrently), so concurrent transactions are those transaction which occurs simultaneously.

Ex- (Ex: Transfer 10\$ from A to B) $\xrightarrow{\text{T1: add log to B}}$

Both are 2 diff-sql answer { T2: Transfer 10\$ from B to A } $\xrightarrow{\text{T1: sub 10$ from B}}$
 $\xrightarrow{\text{T2: add log to A}}$

Mail: gatecess.com appliedroots.com

Consider these 2 sql queries as ~~two different processes in a system if the processes are performed in serial i.e. s_1 , then followed by s_2 , then~~ ^{or vice versa} these will be called as sequential or serial schedule (A schedule is an order of execution of operations).

So $s_1 \rightarrow s_2$ or $s_2 \rightarrow s_1$ are sequential schedule.

Now suppose the transaction sequence is as follows

T_1 : Sub 10\$ from A	}	This is called concurrent or interleaved process of transaction in case of SQL
T_2 : Sub 10\$ from B		
T_1 : Add 10\$ to B		
T_2 : Add 10\$ to A		

*failure
on
crash*

Goal of an OS: Higher throughput (complete as many transactions as possible in given or lesser time)

Solutions

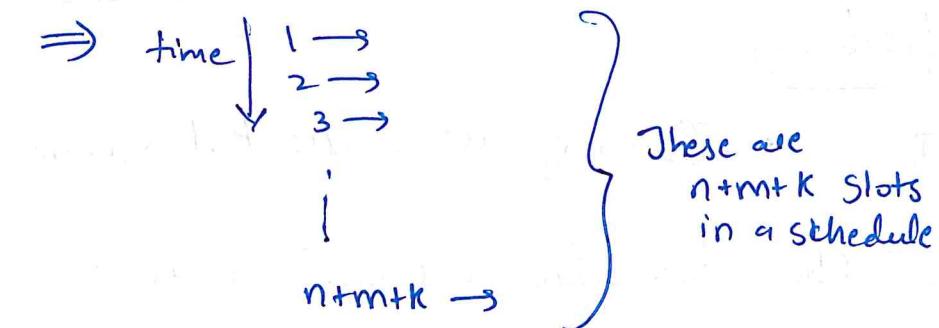
All transactions will be carried out and executed as if it is the only transaction in the system (though there are n no. of transactions taking place simultaneously)

Given 3 transactions T_1, T_2, T_3 with n, m and k operations respectively, number of concurrent schedules will be given as

$T_1: n$ ops $T_1: o_1$
 $T_1: o_2$
 $T_1: o_n$
 $T_2: m$ ops
 $T_3: k$ ops

Phone: +91 844-844-0102

schedules?



So, steps will be

① pick n slots from $n+m+k$ slots

$${}^{n+m+k}C_n$$

now slots left are: $m+k$

② pick m slots from $m+k$ slots

~~$m+k$~~
 ${}^{m+k}C_m$

③ pick ~~the~~ remaining k slots

kC_k

So total no. of schedules are = ${}^{n+m+k}C_n \cdot {}^{m+k}C_m \cdot {}^kC_k$

$$= \frac{(n+m+k)!}{n! m! k!}$$

Total no. of serial schedule = (no. of schedules)!

If no. of transactions given are T then total no. of

Concurrent schedules are = $\frac{(n+m+k)!}{n! m! k!} - T$

Mail: gatecse@appliedroots.com

So, Isolation says, ensure that DB state would be same if each transaction was executed in isolation and sequentially without interleaving. (concurrent schedule execution should be same as any of the serial schedule)

*- Durability:

A DB should be able to recover from any case of failure. OR' Once a transaction is committed its effect must be stored in a non-volatile memory which itself must recoverable from failure.

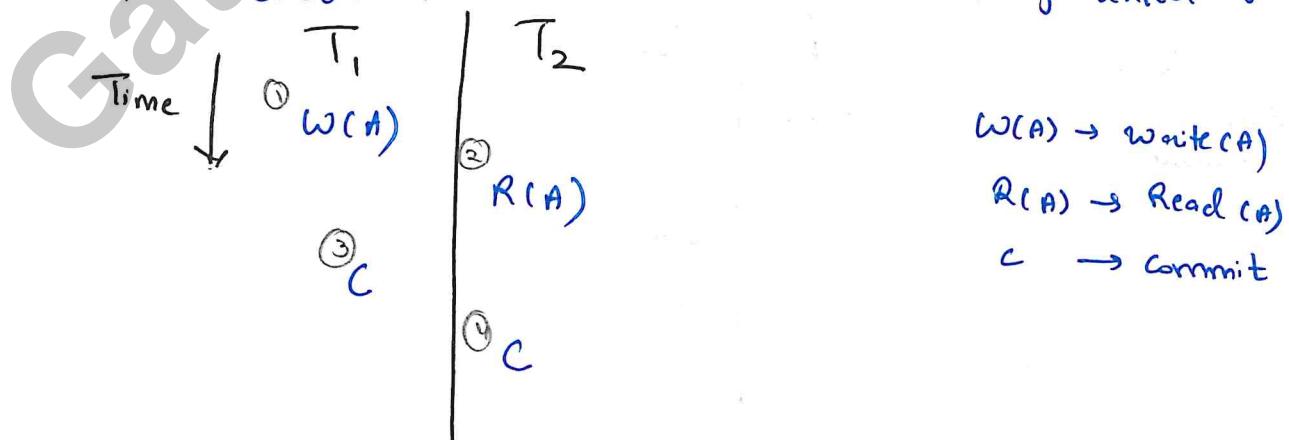
18.3 Problems due to Concurrency:

Objective: higher throughput

Limitation: Hard to Inconsistent Schedule.

*- Write-Read (WR) Problem / Dirty read / Uncommitted read:

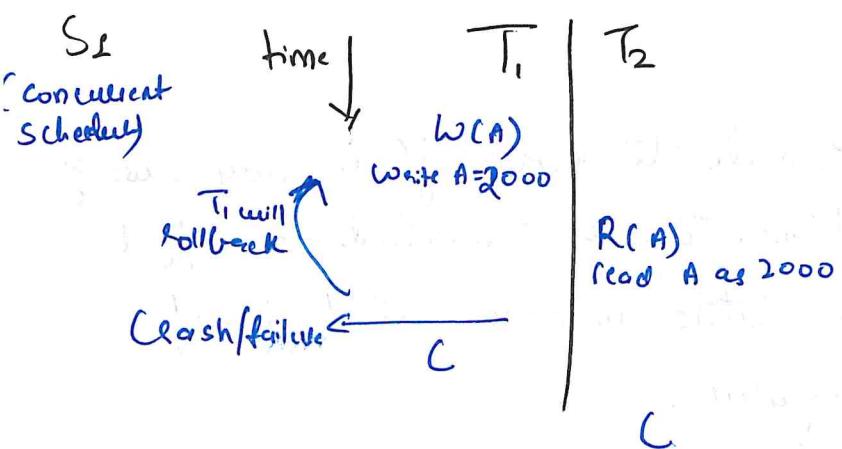
Given a schedule, $R_1(A), R_1(A), C_1, C_2$, represented as:



$w(A) \rightarrow$ write(A)
 $R(A) \rightarrow$ Read(A)
 $C \rightarrow$ Commit

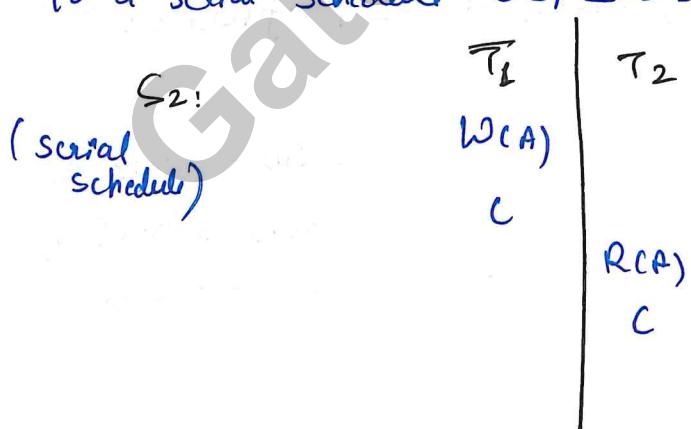
This is called WR problem because we are reading something which is not yet committed.

Suppose, T_1 fails after the execution of T_2 \oplus operation.
So, when T_1 crashes, it will roll back (with the help of log based recovery)



$\therefore T_2$ has already read the initial value and it will continue further with commit operation, this is called dirty read as it is wrong value which T_2 has read (because it ~~read~~ is uncommitted read)

WR-problem may or may not create issues/inconsistency.
Imagine if no crash happened then it will be equivalent to a serial schedule ($S_1 \Leftrightarrow S_2$)

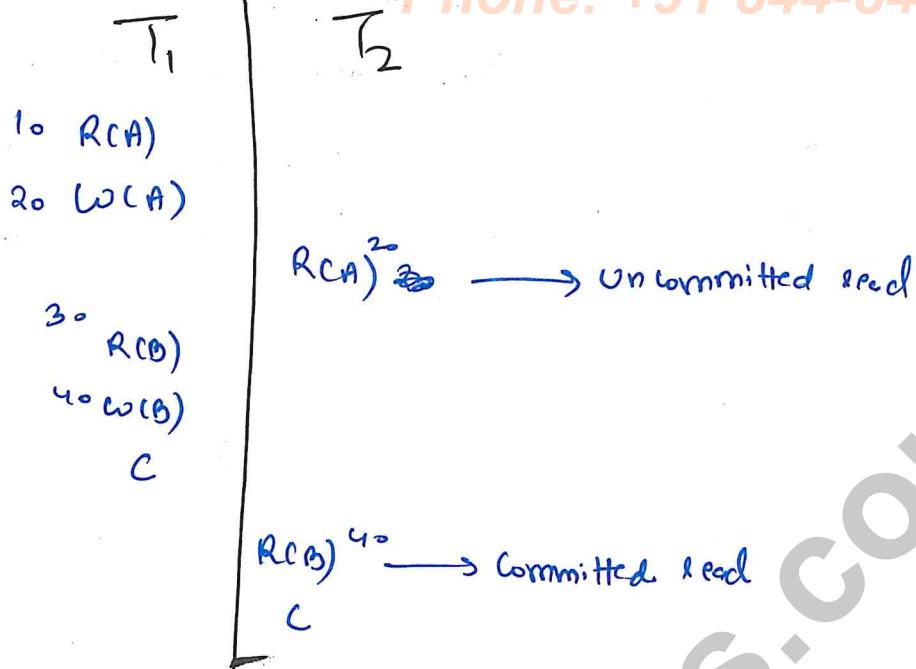


S_1 is said to be serializable or consistent.
 Mail: gatecse@appliedroots.com

Example.

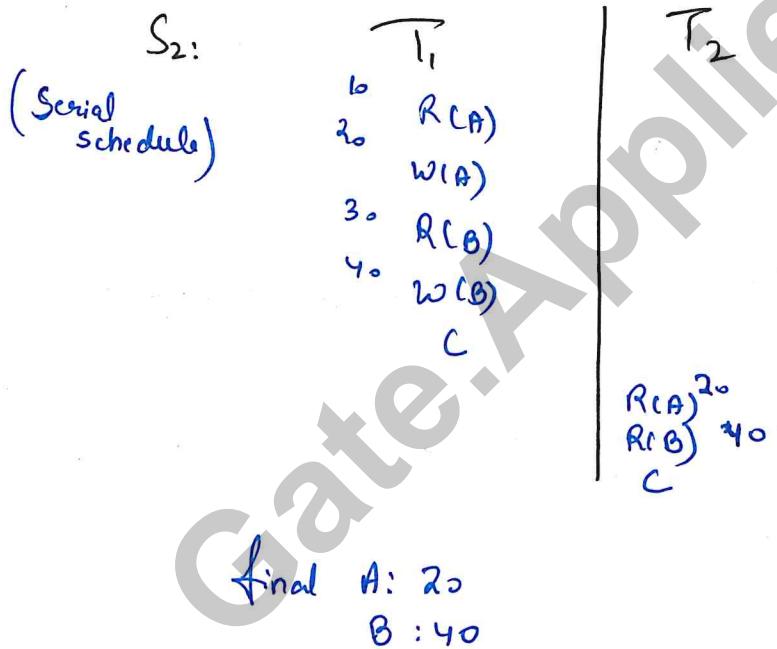
Phone: +91 844-844-0102

APPLIED
ROOTS



If there is no clash, final value of
A: 20
B: 40

In this case it is equivalent to a serial schedule



$$\therefore S_1 \approx S_2$$

hence, S_2 is consistent schedule.

(Here, we have assumed that there is no clash)

Mail: gatecse@appliedroots.com

Note: If commit is not clearly mentioned after all operations are done, there is a default commit.

Example -

S ₁ :	T ₁	T ₂
	R(A) 10	
	W(A) 20	
		R(A) 20
		R(B) 30
		C
	R(B) 30	
	W(B) 40	
	C	

If no crash happen then final values A: 20
B: 40

Let's check whether it is equivalent to serial schedule or not (if no crash happen)

S ₂ :	T ₁	T ₂
	R(A) 10	
	W(A) 20	
	R(B) 30	
	W(B) 40	
	C	

final values A: 20
B: 40 , The values finally we get are same

But the values read by T₂ finally are different.

Hence S₁ \neq S₂. Because even though final values are same but T₂ is reading wrong values.

Let's compare it some other serial schedule **Phone: +91 844-844-0102**

T_1	T_2
	$R(A)$ ¹⁰
	$R(B)$ ³⁰
	C
$R(A)$ ¹⁰	
$w(A)$ ²⁰	
$R(B)$ ³⁰	
$w(B)$ ⁴⁰	
C	

final value of A: 20
B: 40

But still read of T_2 is not same as that of T_1
hence $S_1 \neq S_3$. S_2 and S_3 are only two possible
serial schedule and $S_1 \neq S_2, S_2$

$\therefore S_1$ is inconsistent schedule.

* Read-Write problems (RW) / Unrepeatable reads:

$S_1:$	T_1	T_2
	$R(A)$ ¹⁰	
{ 2 diff values by reading same variable A}		$R(A)$ ¹⁰ $w(A)$ ²⁰ commit
	$R(A)$ ²⁰	
	$w(A)$ ³⁰	
	Commit	

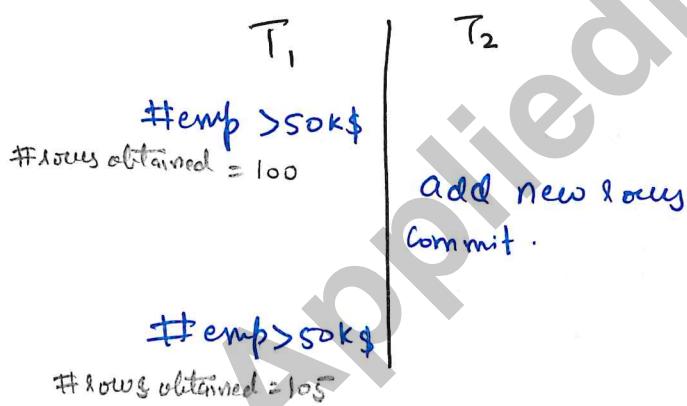
Here, read (A) are performed twice in T_1 and both
gives the diff value due to write operation performed by T_2
This problem is called unrepeatable read problem

Because we read once and then repeated the reads and values are changed as there is read followed by write.

(No consistency in the values that some transaction is reading)

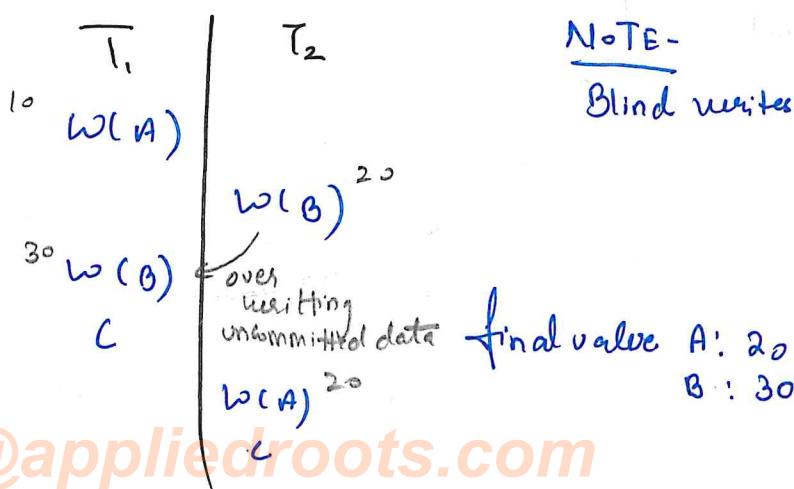
Example- In flight reservation system, suppose 2 person are trying to book the same ticket. Suppose P₁ reads the price and tries to book the ticket while in between P₂ has booked the ticket. As soon as P₁ reaches the paid payment option, it is showing no tickets are available.

There is a special case of read-write problem called phantom read problem. In this new rows are added or deleted in between.



→ Write-Write Problem (WW) / overwriting uncommitted data:

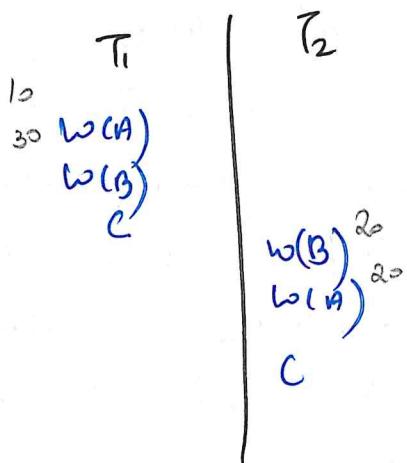
S1:



Let S_2 be a serial schedule: (Suppose that no clash happened in S_1)
Phone +91 844-844-0102

APPLIED
ROOTS

S_2 :

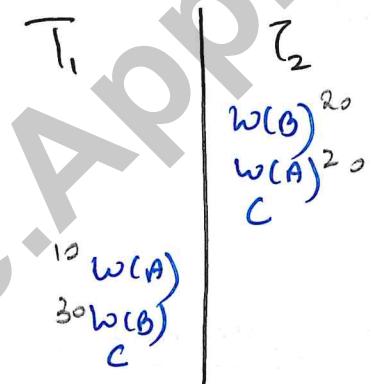


final values of A: 20
B: 20

Here, final values of S_1 and S_2 do not match. Hence
 $S_1 \neq S_2$.

Let S_3 be another serial schedule:

S_3 :



final values of A: 10
B: 30

So, $S_1 \neq S_2 \neq S_3$.

$\therefore S_1$ is inconsistent schedule.

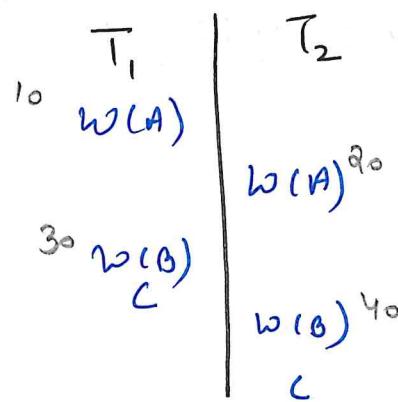
Note that W-W problem may or maynot create inconsistency

Mail: gatecse@appliedroots.com

Example

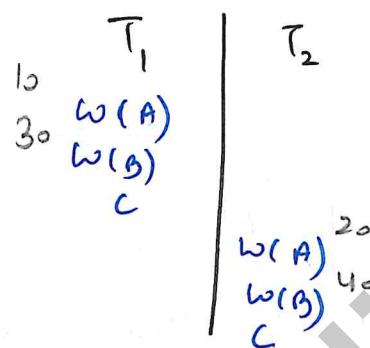
APPLIED
ROOTS

$S_1:$



final value of A: 20
B: 40

$S_2:$



final value of A: 20
B: 40

So, here $S_1 \approx S_2$. So WW may or may not create inconsistency.

Note If a schedule say S_1 is consistent than there is no problem. Problem arises when we have concurrent schedule which is inconsistent or schedule crashes.

Case 1: When there is clash in w-w conflict **Phone: 91 844-844-0102**
problem is called as lost-update (WW-rollback)

S_t:

	T ₁	T ₂
10	R(A)	
20	w(A)	
		w(A) ³⁰
	c	c

→ lost the changes made by T₂

crash

logs:

T₁: A:10 → 20

T₂: A: 20 → 30

T₂: C

~~T₁~~ crash ($\because T_1$ has not committed yet \therefore we will roll back T_1)

If T_1 is rolled back then all the updates made by T_2 will get lost (this lost the changes made by T_2 even though T_2 didn't crashed)

This problem is called lost update problem.

Note - WW-issue + inconsistent \rightarrow lost update in case of crash + rollback.

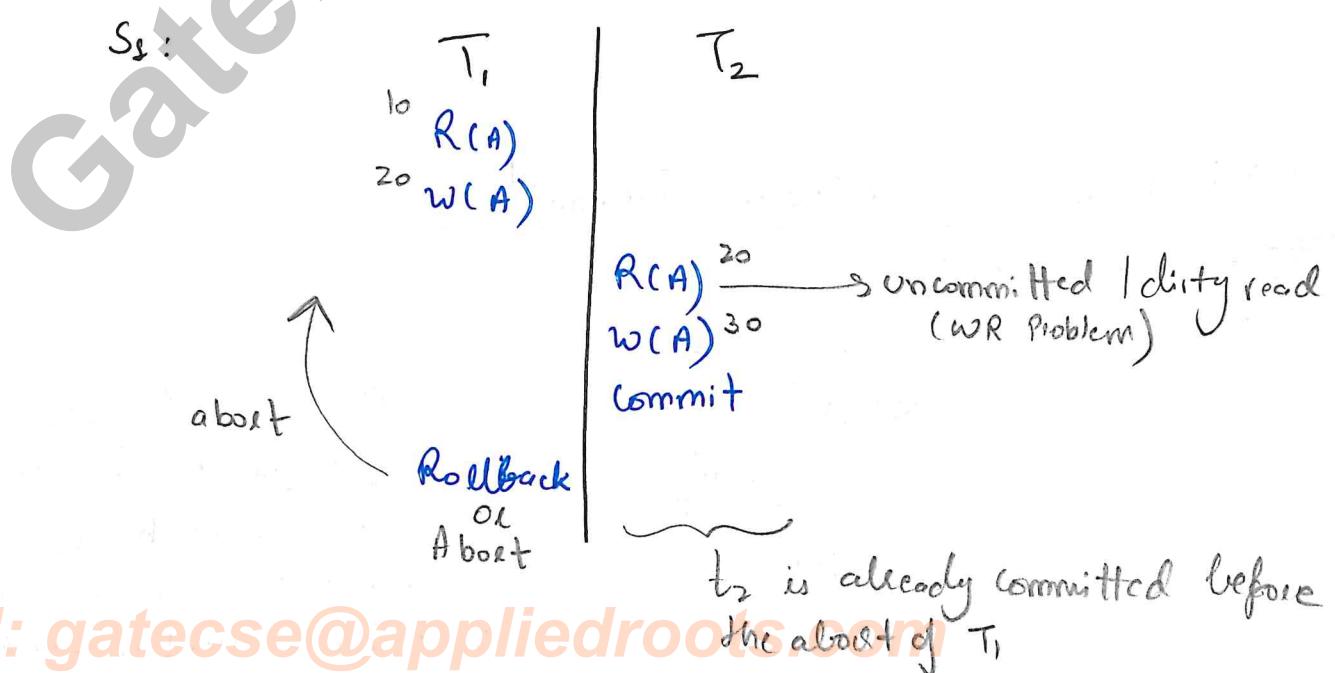
For the given concurrent schedules, we have 2 important properties that they should satisfy:

- Recoverability
- Serializability

Recoverability states that: If a transaction is Aborted (ended unexpectedly) or some crash happened and the transaction is rolled back then we should be able to reach to the initial state. ie the schedule should be recoverable. (We want our Schedule to be crash proof)

Serializability is concerned to isolation. When a concurrent schedule is equivalent to serial schedule then we say that concurrent schedule is serializable.

→ Irrecoverable Schedules:

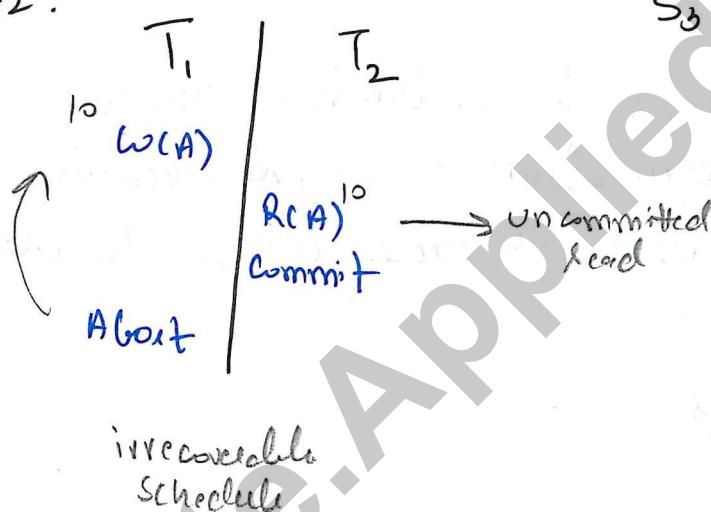


This schedule S_2 is not recoverable because T_2 has committed the values and now it cannot be undone. So the decisions made here are ~~can~~ cannot be undone as it has been committed. If the commit of T_2 is placed after the abort of T_1 , then there would be no issues.

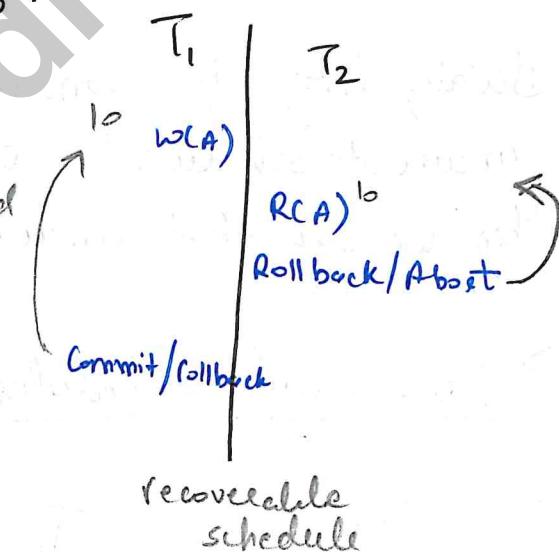
So, Rolling back a committed transaction is irrecoverable.

Example -

$S_2:$



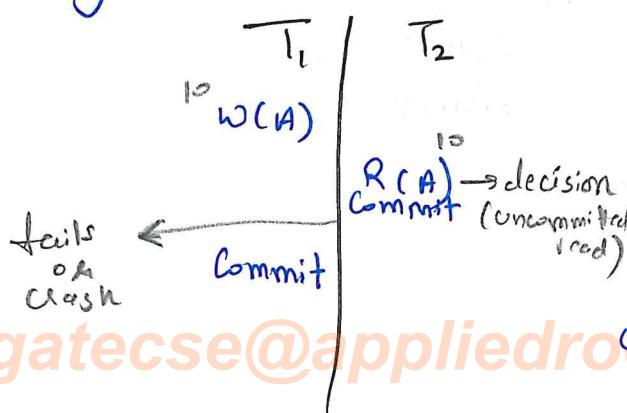
$S_3:$



→ Dependence between Transaction:

(case 1:
 $S_1:$

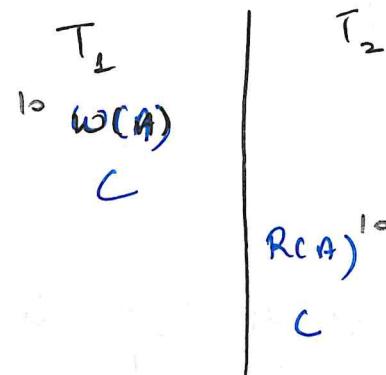
initially $A=5$



On this case T_2 depends on T_1 because decision made by T_2 is either valid or invalid based on whether T_1 will rollback or not.

Phone: +91 844-844-0102

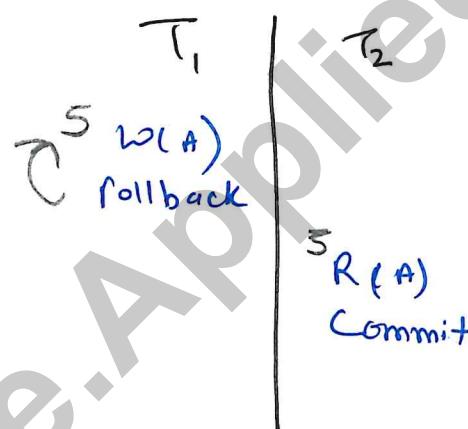
So, in S_1 validity of transaction 2 (T_2) decision is based on T_1 success and failure. That is why T_2 depends on T_1 .

Case 2: $S_2:$ 

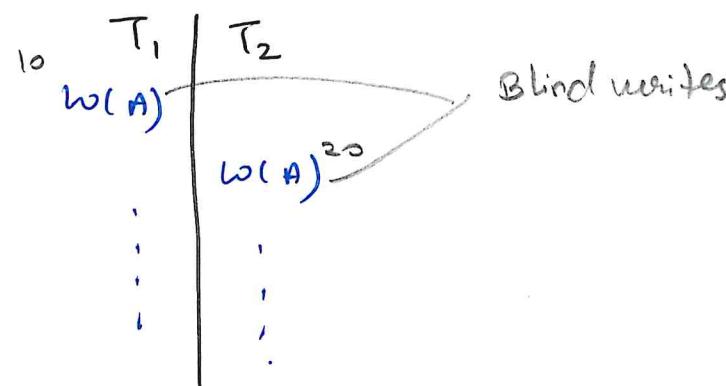
T_2 has read the committed value.

The failure or validity of value in T_2 doesn't depend on T_1 in S_2 because T_1 has already committed.

In this case T_2 does not depend on T_1 .

Case 3: $S_3:$ 

here also T_2 does not depend on T_1

Case 4: $S_4:$ 

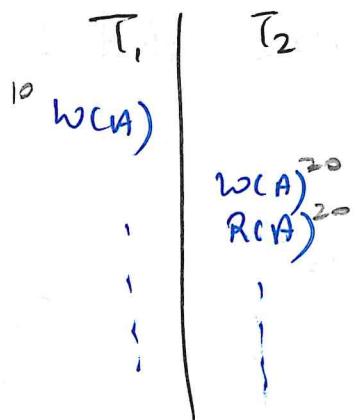
Mail: gatcse@appliedroots.com

here, T_2 is not dependent on T_1 as T_2 is not reading any value written by T_1 as no actual change is made by T_1 .

Case 5 :

Phone: +91 844-844-0102

S5:
APPLIED
ROOTS

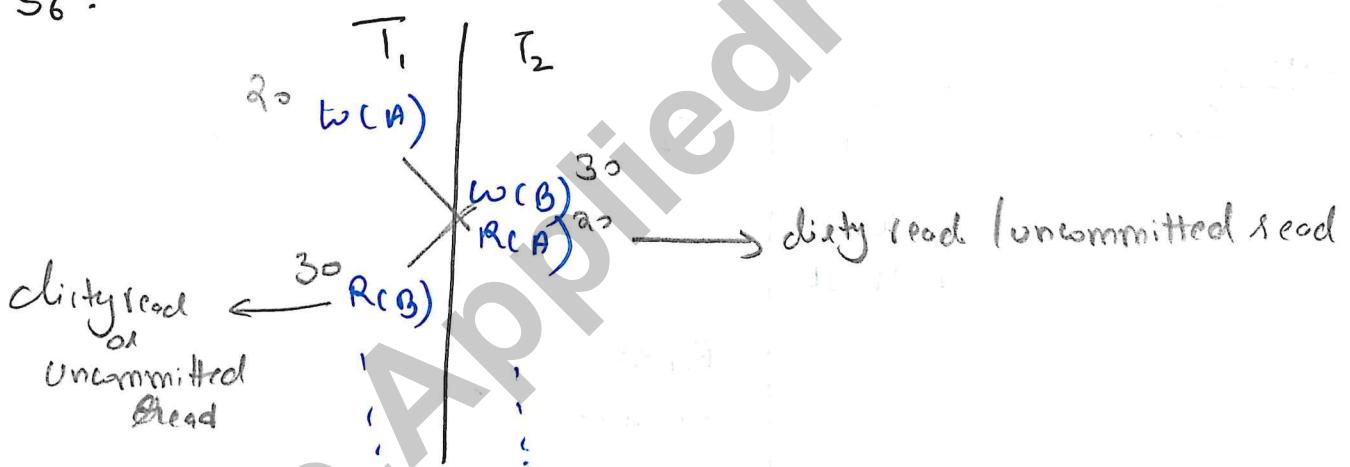


→ This is not an uncommitted read because what we are reading here is written in the same transaction.

T₂ doesn't depend on T₁ as R₂(A) is not a uncommitted read, it has read the value written in T₂ itself.

Case 6 :

S6 :



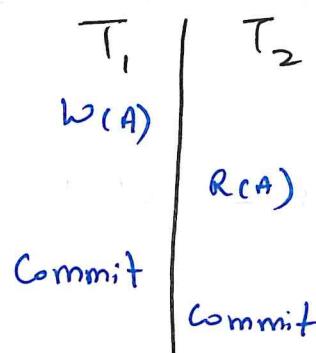
T₂ is dependent on T₁ and T₁ is also dependent on T₂.

Mail: gatecse@appliedroots.com

→ How to ensure recoverability of schedules? **Phone: +91 914844-0102**

If T_2 depends on T_1 then T_2 should commit only after T_1 commits or rollback. If this condition is satisfied that it is guaranteed to have a recoverable schedule.

$S_1:$

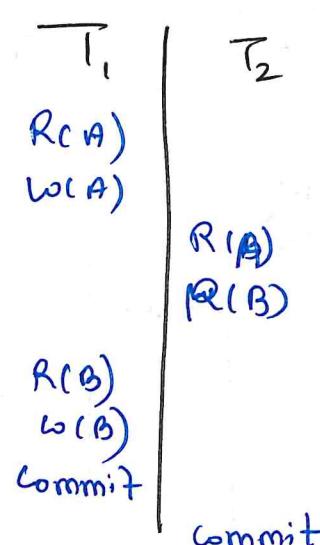


Here, T_2 is dependent on T_1 and T_2 has committed after the T_1 has committed. Therefore, S_1 is recoverable schedule.

Note: A recoverable schedule may have WR (dirty read) problem.

Example:

$S_2:$



Is S_2 recoverable?

here T_2 depends on T_1

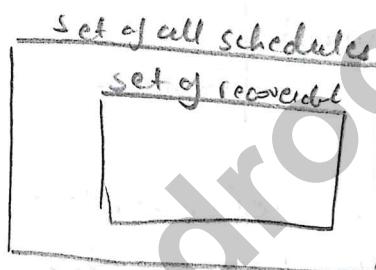


S_2 is recoverable iff
 C_2 (commit by T_2)
happens after C_1

$\therefore T_2$ committed after $T_1 \therefore S_2$ is recoverable.

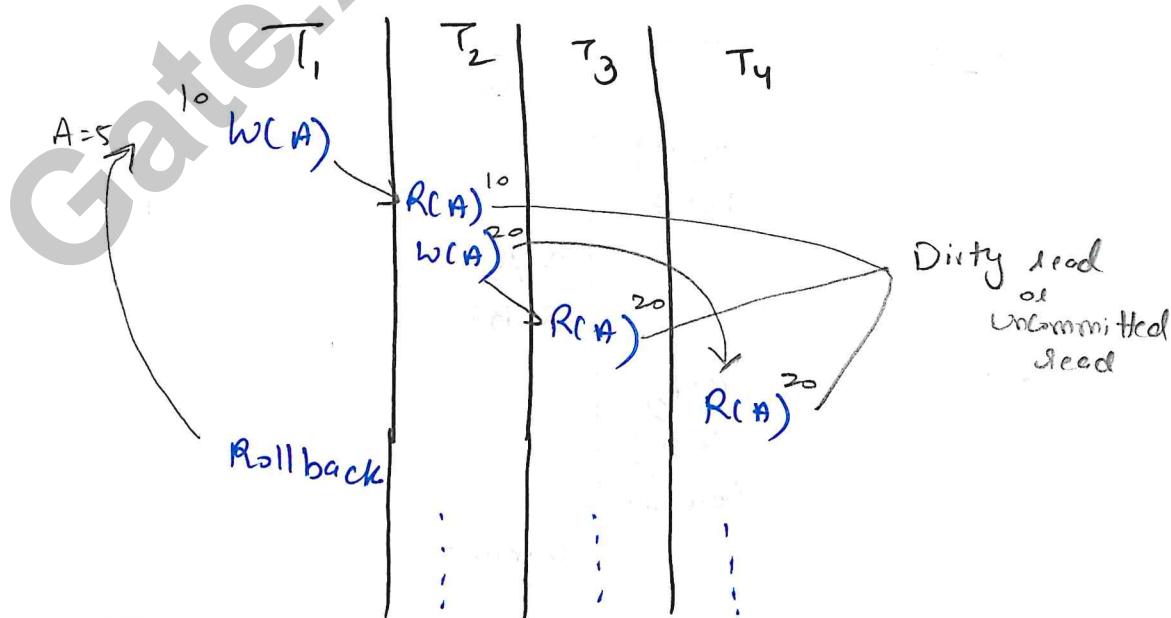
NOTE -

Set of all schedules \supseteq set of recoverable schedules



Cascading Rollback

S_1 : initially $A = 5$



Mail: gatesschedule@appliedroots.com The schedule S_1 is recoverable as all rollbacks are ordered as no schedule has committed

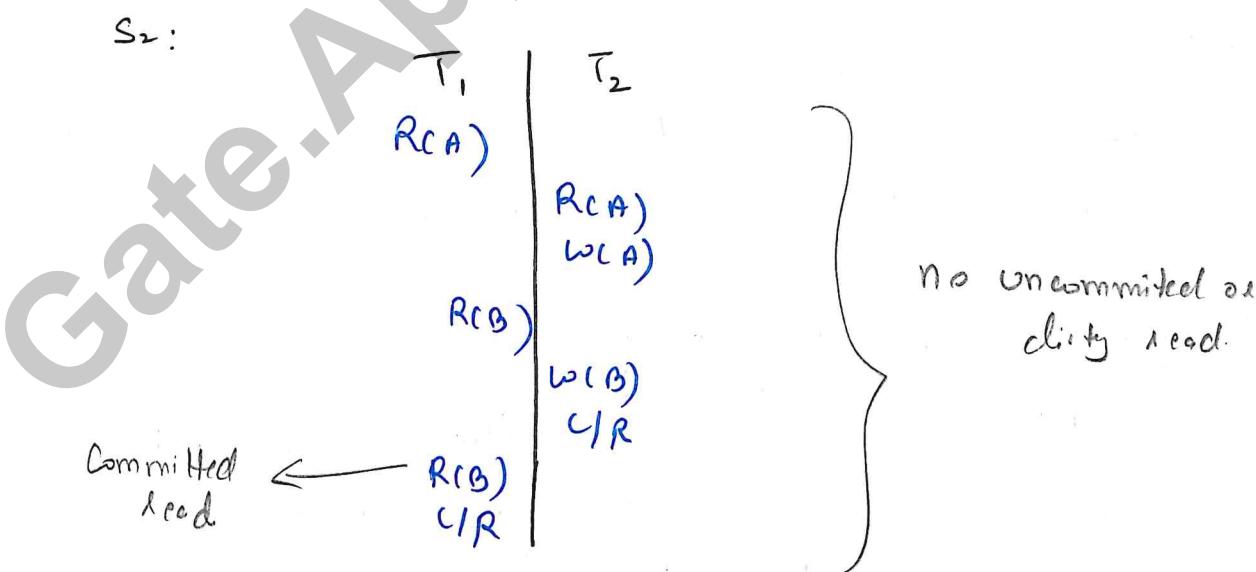
But, here rollback of T_1 will result in rollback of T_2 . Rollback of T_2 will result in rollback of T_3 and so on. (cascading rollbacks). It is recoverable but it causes other problem like lower throughput and longer wait time.

T_2 , T_3 , T_4 doesn't any mistake but they are dependent on T_1 \therefore they have to rollback.

We should avoid the cascading rollbacks.

→ How to avoid Cascading Aborts (ACA) / cascadeless:

Cascading aborts are happening primarily because of dirty reads. So, in order to avoid that there should not be presence of uncommitted reads. i.e disallow a transaction from reading uncommitted changes from another transaction.



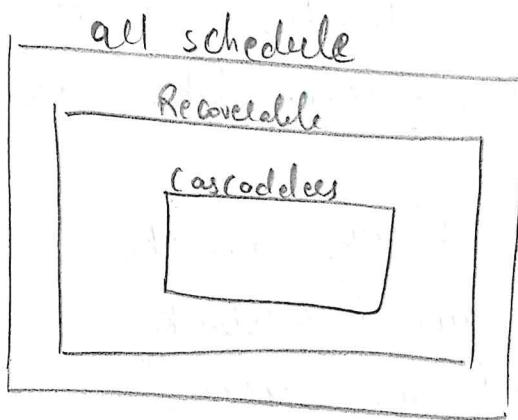
S_2 is both cascadeless and recoverable schedule.

NOTE -

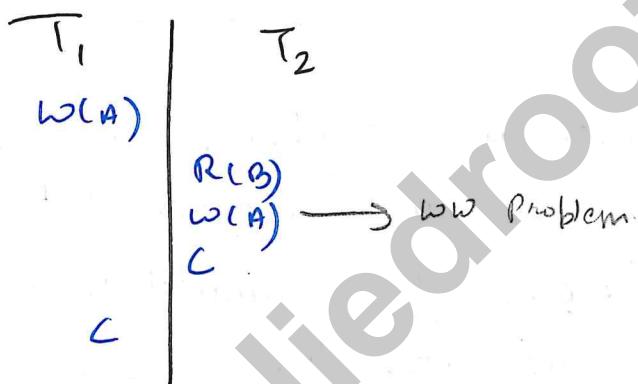
Phone: +91 844-844-0102

APPLIED
ROOTS

Set of recoverable schedule \supseteq Set of cascades schedule



S3:



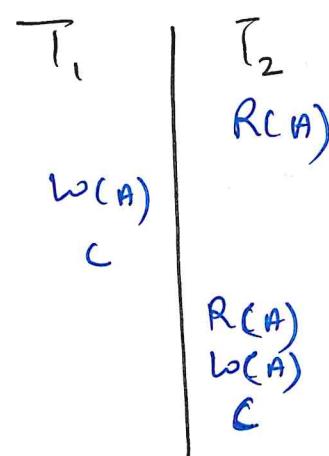
S_3 is cascades and recoverable but it still have WW problem.

Cascades guarantees that there will be no WR problem
or dirty problems, other problems may exist.

→ Strict schedules

If T_1 performs $w(A)$ then T_2 is not allowed to $R(A)$ or $w(A)$ till T_1 commits or rollbacks. (ie No WR and No WW)

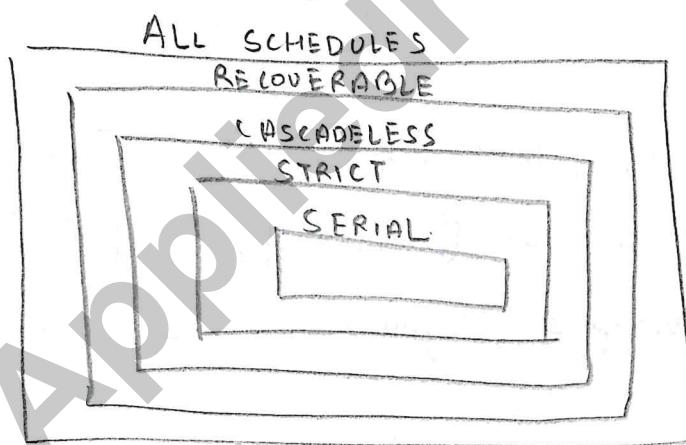
Mail: gatecse@appliedroots.com

$S_1:$ 

S_1 is strict schedule (it is both cascadeless as well as recoverable)

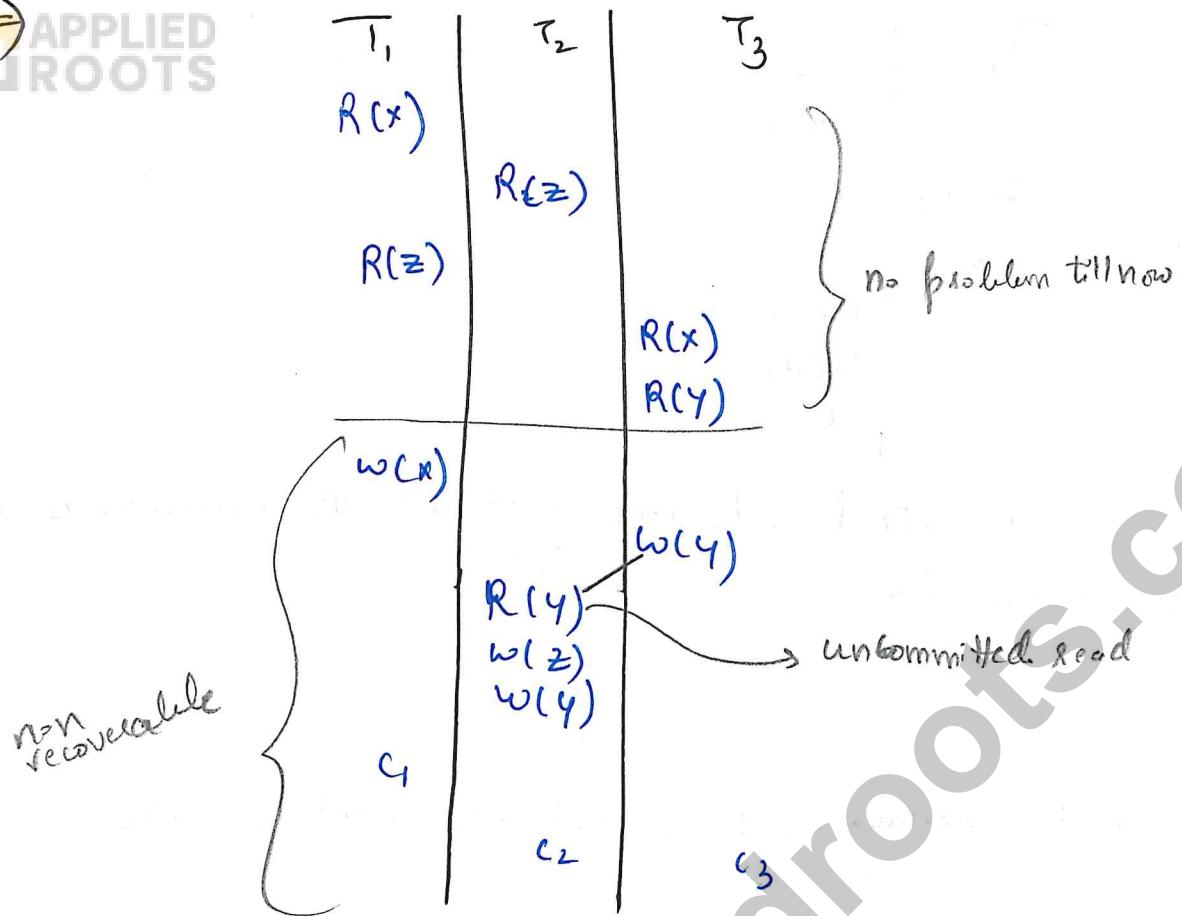
Note -

Complete ~~legal~~ set of schedules will be like:



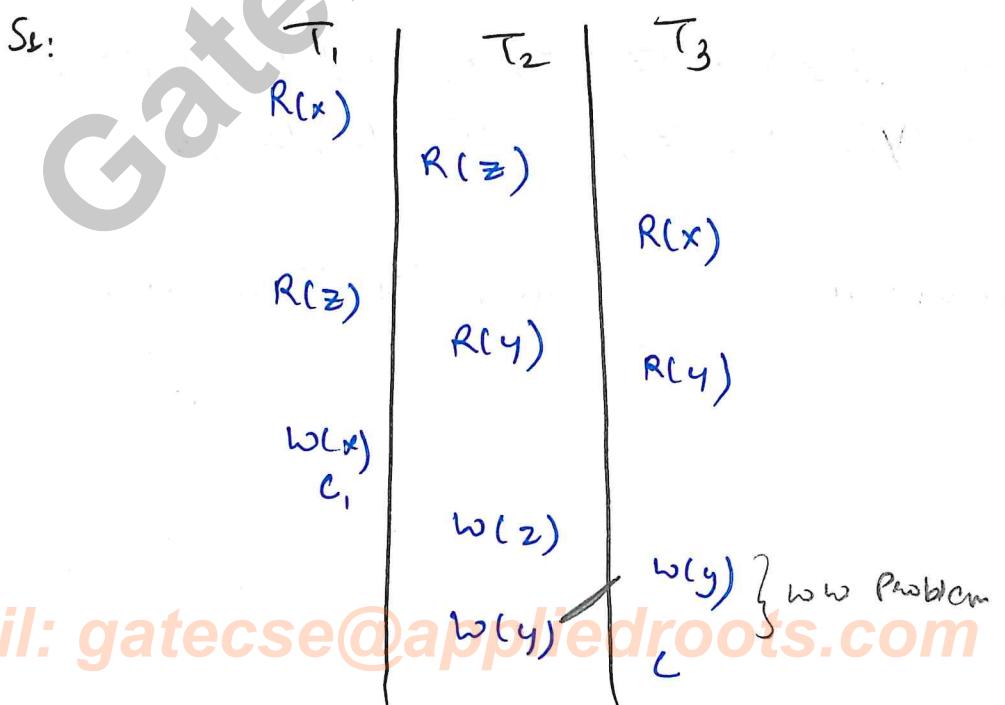
(Q1) $S_1: R_1(x), R_2(z), R_1(z), R_3(x), R_3(y), w_1(x), w_3(y), R_2(y), w_2(z), w_2(y), C_1, C_2, C_3$.

S_1 is the schedule along with its time order.



Due to uncommitted read in T_2 , T_2 depends T_3 . So S_2 to be recoverable if C_2 happens after C_3 . But since C_2 happens before C_3 in S_2 $\therefore S_2$ is not recoverable

(Q2) Qs S_2 is recoverable, cascades, strict or not?



⇒ Since, there are no dependency $\therefore S_2$ is recoverable (no dirty reads) **Phone: +91 844-844-0102**
but there are LRU problem.

\therefore No dependency $\rightarrow S_2$ is recoverable

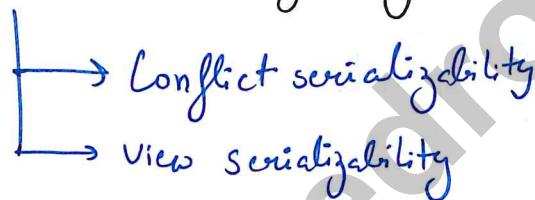
No dirty reads $\rightarrow S_1$ is cascades

\therefore we have LRU issue \rightarrow NOT STRICT.

18.5 Serializability of schedules - Conflict serializability:

A concurrent schedule to be serializable should be equivalent to any of the serial schedule.

→ How to determine serializability?

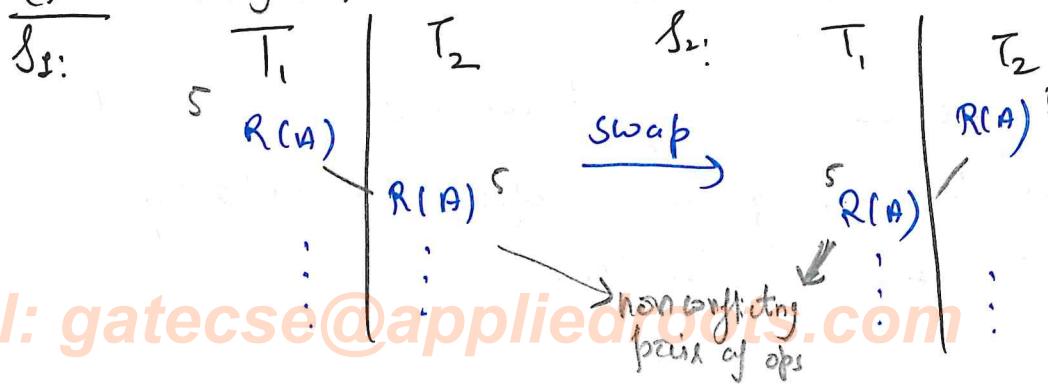


Before getting into conflict serializability we need to get to know some/few terminologies.

* Non-conflicting pair of operations:

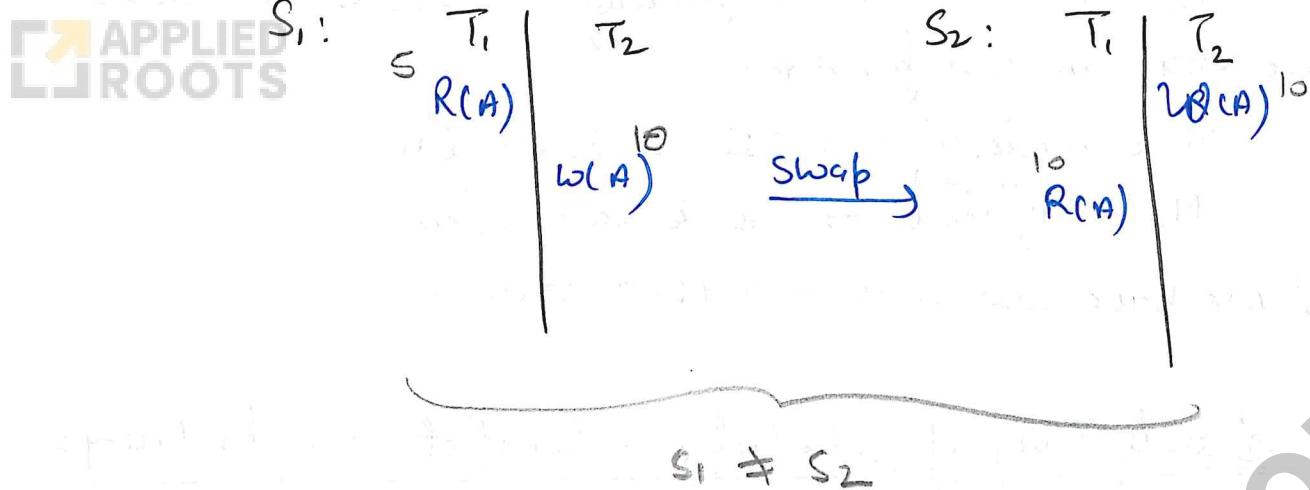
On swapping a pair of operations, if the resulting schedule is equivalent to the original schedule than those pair of operation are called non-conflicting pair of operations.

Ex- initially $A=5$



$S_1 \approx S_2$
such pair of ops are called non-conflicting operations

Ex-2 initially $A = 5$



$\therefore S_1 \neq S_2 \therefore$ the pair of ops given is a conflicting pair of operations.

→ Conflicting pair of operation:

A pair of operation is conflicting, if:

- * Both operations should not belong to the same transaction
 - * Both operations must be on the same data
 - * At least one of the operation should be "write"
- If ~~any~~ these 3 conditions are satisfied than the resultant pair of operation is a conflicting pair of operation.

→ Conflict-equivalent schedules:

A schedule S₁ is conflict equivalent to S₂ if

S₁ results in S₂ from swapping non conflicting pairs. Let S₁ and S₂ be:

S_1 APPLIED
ROOTS T_1 $R(A)$ $w(A)$ non conflicting
pair of ops $R(B)$
 $w(B)$ T_2 Conflicting
pair of ops $R(A)$ $R(B)$

:

 \approx_c T_1 $R(A)$ $w(A)$ $R(B)$ $w(B)$ $R(A)$ $R(B)$

here we have
swapped conflicting
pair of ops
 $w_1(A)$ and $R_2(A)$

Q1 Determine whether S_1 and S_2 are conflict equivalent or not.

$\Rightarrow S_1:$

T_1	T_2	T_3
	$R(A)$ $w(A)$	
	$w(B)$	$R(C)$
$R(A)$ $R(B)$ $w(A)$ $w(B)$		$w(A)$ $w(C)$
:	:	:

 $S_2:$

T_1	T_2	T_3
	$R(A)$	
	$w(A)$ $w(B)$	$R(C)$ $w(C)$
$R(A)$ $R(B)$ $w(A)$ $w(B)$		$w(C)$
:	:	:

S_1 and S_2 are not conflict equivalent because of $w_2(A)$ and $w_3(A)$

→ Conflict Serializable schedule: (LSS) **Phone: +91 844-844-0102**

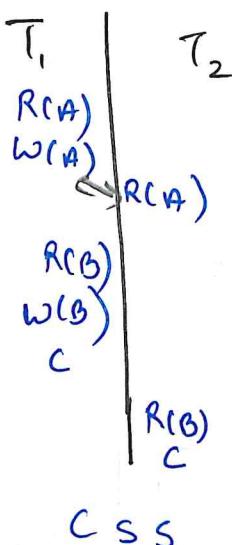


A schedule s is CSS if there exist a serial schedule s' such that s is conflict equivalent to s'

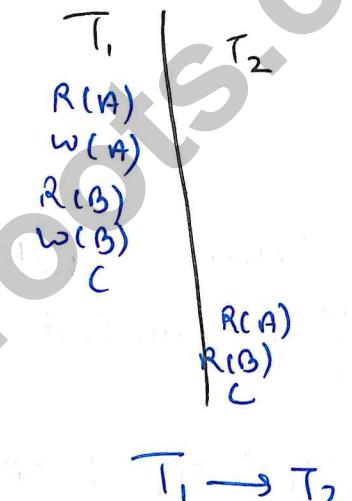
i.e. $s \approx s'$

Concurrent schedule \leftarrow $s \approx s'$ \rightarrow serial schedule

$s:$



$s':$



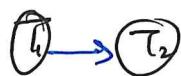
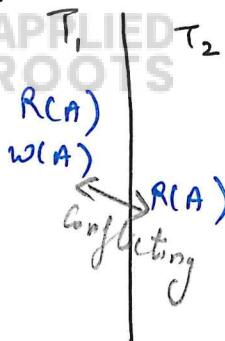
$T_1 \rightarrow T_2$

We have swapped conflicting pair

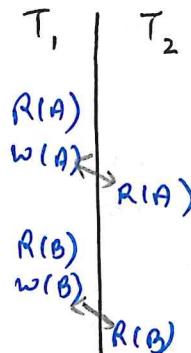
Simpler Method to determine CSS is Precedence Graph.

The steps to draw precedence graph is:

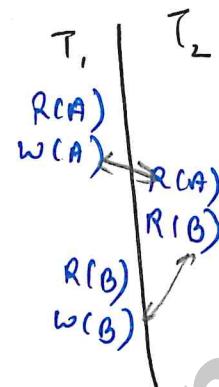
- [1] Every transaction is a vertex
- [2] Insert/Create a directed edge between $T_i \rightarrow T_j$ if \exists a conflicting pair and T_i proceeds T_j

S₁:

T_1 occurs before T_2

S₂:

T_1 occurs before T_2

S₃:

T_1 forms a cycle because T_1 and T_2 are dependent on each other.

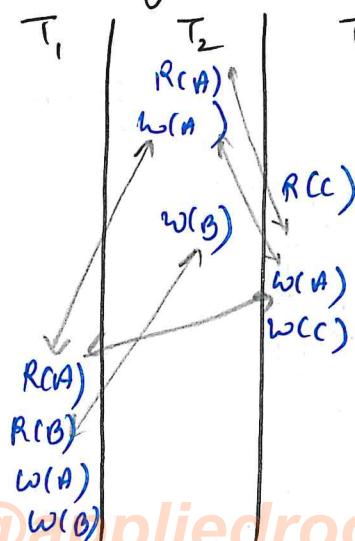
With the help of the precedence graph we can check whether a schedule is CSS or not:

If there is a cycle in precedence graph than that schedule is not CSS while vice versa is true. i.e if there is no cycle for a given schedule than it is CSS.

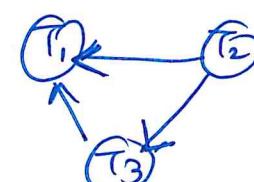
$\therefore S_1$ and S_2 are CSS while S_3 is not CSS.

Using precedence graph we can easily determine whether given schedule is CSS or not.

(Q1) Determine whether given schedule is CSS or not?



So, T_3 and T_1 depends on T_2 and T_1 depends on T_3 as well

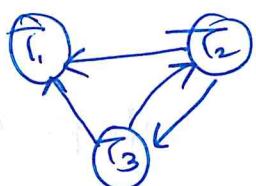
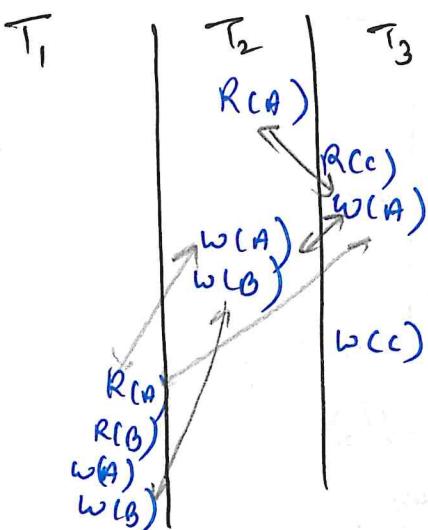


Hence, it is CSS as there is no cycle

(Q2) Is it CSS or not?

Phone: +91 844-844-0102

APPLIED
ROOTS



∴ there is cycle ∴ it is not CSS.

Dr. Pradeep's question

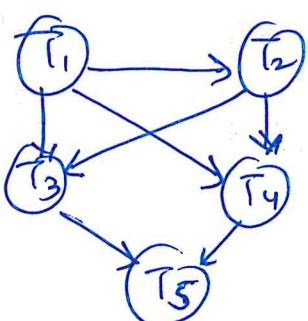
Trace

→ How to get equivalent serial schedule: (To topological sort)

- (1) Start with T_i whose indegree is zero.
- (2) Remove T_i and all edges connecting from T_i
- (3) Repeat 2 and 1

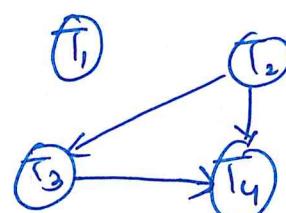
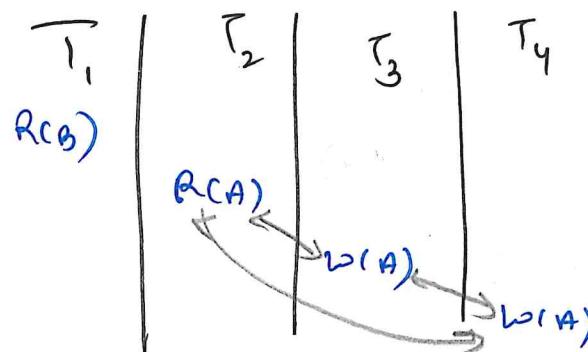
It gives the serial schedule.

Ex- If the graph (Precedence graph) is given as



On applying topological sorting 2 serial schedules we will obtain are:- $T_1 \rightarrow T_2 \rightarrow T_3 \rightarrow T_4 \rightarrow T_5$, $T_1 \rightarrow T_2 \rightarrow T_4 \rightarrow T_3 \rightarrow T_5$

(Q3) Is it CSS or not?



\therefore there is no cycle

\therefore it is CSS

It is equivalent to serial schedule $T_2 \rightarrow T_3 \rightarrow T_4$. ~~From~~

T_1 could be placed ~~anywhere~~ anywhere

\therefore we have 4 serial schedule ~~which~~ to which this schedule is equivalent to.

$$T_1 \rightarrow T_2 \rightarrow T_3 \rightarrow T_4$$

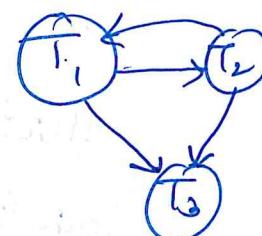
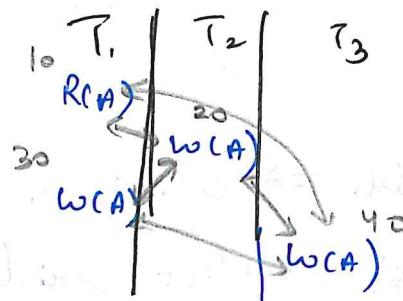
$$T_2 \rightarrow T_1 \rightarrow T_3 \rightarrow T_4$$

$$T_2 \rightarrow T_3 \rightarrow T_1 \rightarrow T_4$$

$$T_2 \rightarrow T_3 \rightarrow T_4 \rightarrow T_1$$

(Q4) Is it CSS?

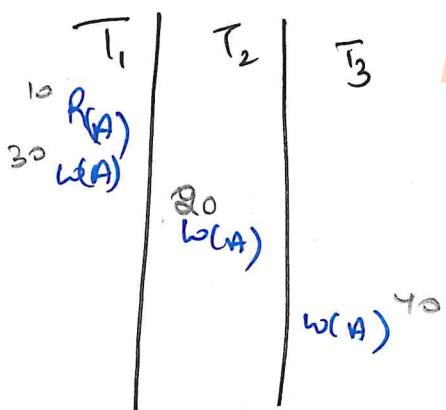
S₃:



\therefore there is a cycle \therefore it is CSS

and hence, it is not equivalent to any serial schedule.

S_2 :



Phone: +91 844-844-0102

$S_1 \not\approx S_2$

but $S_1 \approx S_2$ because for every read operation value read is the same and final value written is also same. So, S_1 is ~~not~~ equivalent to S_2 (S_2 is serial schedule) but S_1 is not conflict equivalent to S_2 .

NOTE -

If schedule is conflict serializable then it is serializable (i.e. it is having an equivalent serial schedule). But every serial schedule / serializable schedule may not be conflict serializable.

The schedules S_1 and S_2 are not conflict serializable though S_2 is serial schedule.

So, CSS is a sufficient condition for serializability NOT a necessary condition.

(If precedence graph is acyclic then CSS \Rightarrow serializable else, NOT CSS but may/may not be serializable.)

Conflict serializable \Rightarrow serializable

Serializable $\not\Rightarrow$ conflict serializable

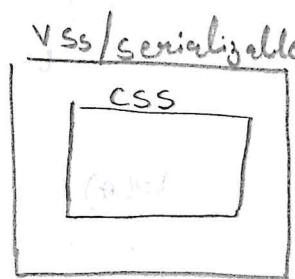
18.6 View Serializability & (VSS)

Phone: +91 844-844-0102

If a schedule is view serializable then the schedule is serializable else it is not serializable.

View serializable \Leftrightarrow Serializable

VSS is both necessary and sufficient condition for serializability.



→ view equivalence:

(Represented by \approx_v)

$S_1 \approx_v S_2$ only if,

[1] Initial reads must match:

Ex - initially $A=10$ (initially) $B=20$

$S_1:$	T_1	T_2	$S_2:$	T_1	T_2
	$R(A)$			$R(A)$	
				$R(A)$	

(Initial reads are same)

Ex - initially $A=10$
 $B=20$

$S_1:$	T_1	T_2	T_3
	$R(A)$		
	$R(B)$		

$S_2:$	T_1	T_2	T_3
		$w(B)$	
		$R(B)$	

here, initial

reads do not

match, as $R(B)$ in S_2
read diff from $R(B)$ in S_1

Mail: gate@appliedroots.com

[2] Write - Read (WR) must match.

Phone: +91 844-844-0102

APPLIED ROOTS

Ex- initially $A = 10$

$S_1:$	T_1	T_2	T_3
	20 $w(A)$		
		30 $w(A)$	

$S_2:$	T_1	T_2	T_3
	20 $w(A)$		
		30 $w(A)$	20 $R(A)$

$S_3:$	T_1	T_2	T_3
	20 $w(A)$		
		30 $w(A)$	30 $R(A)$

In S_1 , $R_3(A)$ read 30 while in S_2 , $R_3(A)$ read 20
 $\therefore S_1 \not\approx S_2$

In S_1 and S_3 , the WR operation ($w_2(A) - R_3(A)$) is not changed \because both the $R_3(A)$ has read some value 30
 $\therefore S_1 \approx S_3$

In S_2 and S_3 , the WR operation is changed
 $\therefore S_2 \not\approx S_3$

[3] Final writers must match for each variable.

Initially, $A = 10, B = 20$

<u>Ex</u>	$S_1:$	T_1	T_2	T_3
		30 $w(A)$		
			40 $w(B)$	

$S_2:$	T_1	T_2	T_3
	30 $w(A)$		
		40 $w(B)$	50 $w(A)$

$\not\approx$

$\not\approx$

$\not\approx$

$\not\approx$

$\not\approx$

Final $A = 50$
 $B = 70$

Mail: gatecse@appliedroots.com

Since, final writes of A and B is not same for S_1 and S_2 \therefore they are not view equivalent.
Suppose we have S_3 as:

$S_3:$	T_1	T_2	T_3
	30 $w(A)$ $w(B)$	50 $w(A)$	
			60 $w(A)$ $w(B)$

here, final write is done by T_3 (same as S_1).
 $A=70$ and same for B \therefore final writes are same
hence $S_1 \approx_v S_3$

So, a schedule S_1 is view equivalent to S_2 if all the three conditions are matched.

→ View Serializability:

If given a schedule S is view equivalent to a serial schedule than it is called "View Serializable Schedule (VSS)".

Ex -

$S_1:$	T_1	T_2	T_3
	$w(A)$		
		$w(A)$	
			$w(A)$

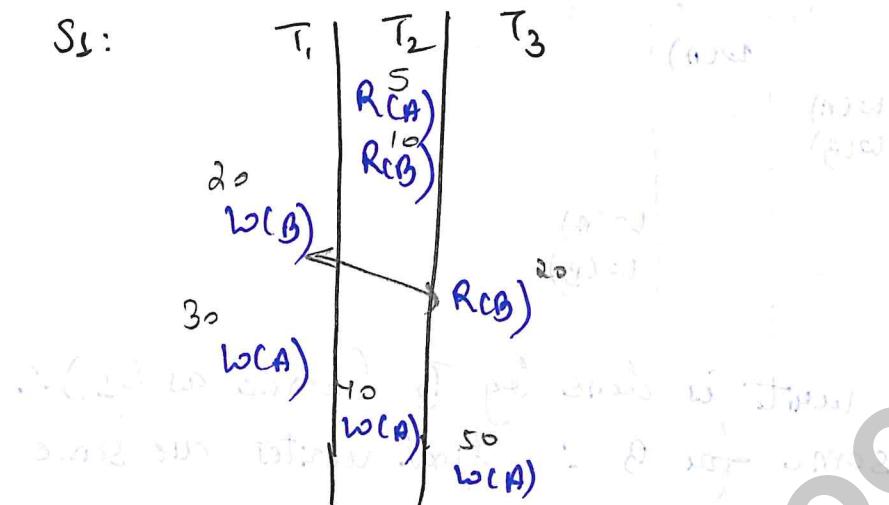
$S_2:$	T_1	T_2	T_3
			$w(A)$
		$w(A)$	
			$w(A)$

\approx_v

Mail: Njatecs.No@appliedroots.com

in both S_1 and S_2 \therefore S_1 & S_2 is VSS but not SSS.

(Q1) Is it VSS?



→ Let initial value of $A = 5, B = 10$

① Initial read: T_2 should be before T_1 ($T_2 \rightarrow T_1$)
 and T_2 should be before T_3 ($T_2 \rightarrow T_3$)

(Because in both cases read will not satisfy)

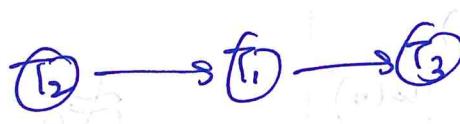
② W-R seq. should not be altered: So $T_1 \rightarrow T_3$

③ Final writes: $T_2 \rightarrow T_3$
 $T_1 \rightarrow T_3$

So T_3 should write last

So from ① ② and ③ serial schedule

will be



∴ S_1 is VSS

The precedence graph for this is:-

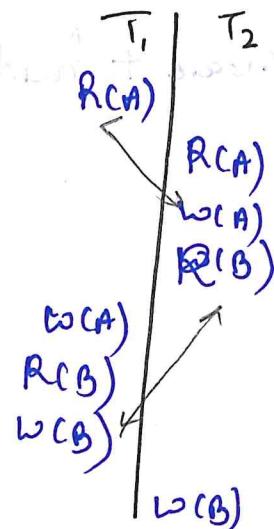


Mail: gatecse@appliedroots.com

∴ There is cycle
 and for this reason S_1 is not CSS.

(Q2) Is it VSS?

S_1 :



→ Let's mark regions:

① Initial read \Rightarrow T_1 should be before T_2 as there is $R_1(A) \rightarrow W_2(A)$
 $\therefore T_1 \rightarrow T_2$

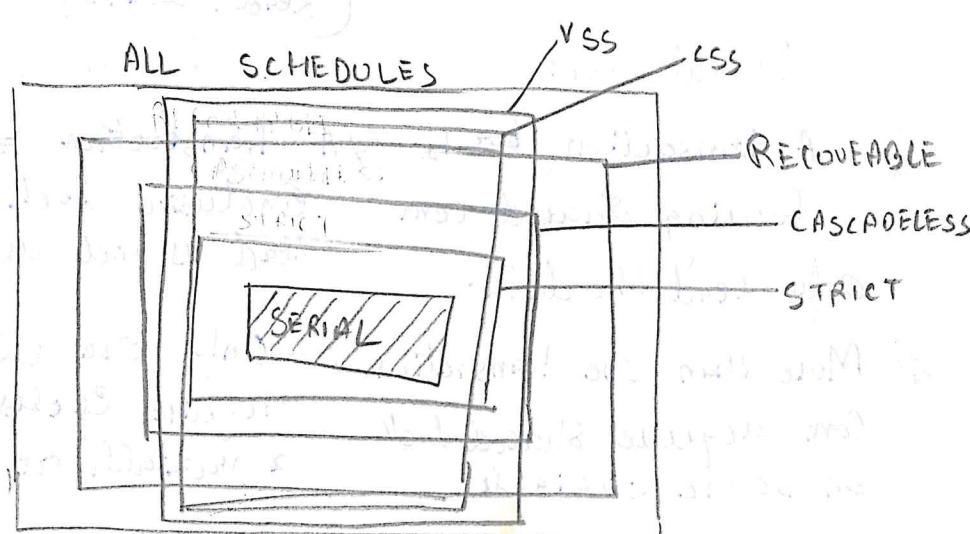
T_2 should be before T_1 as there is $R_2(B) \rightarrow W_1(B)$
 $\text{ie } T_2 \rightarrow T_1$

\therefore It is a conflict
 $\therefore S_1$ is not VSS

$\therefore S_1$ is not VSS \therefore it is not CSS

NOTE -

Recoverability doesn't guarantee serializability.
 Serializability doesn't guarantee recoverability.



18.7 Lock based Concurrency Control

~~Recoverability~~ ~~Consistency~~ ~~Isolation~~ ~~Durability~~ ~~ACID~~ ~~844-844-0102~~

Serializability

Lock is intuitively a software + hardware functionality provided implemented by OS.

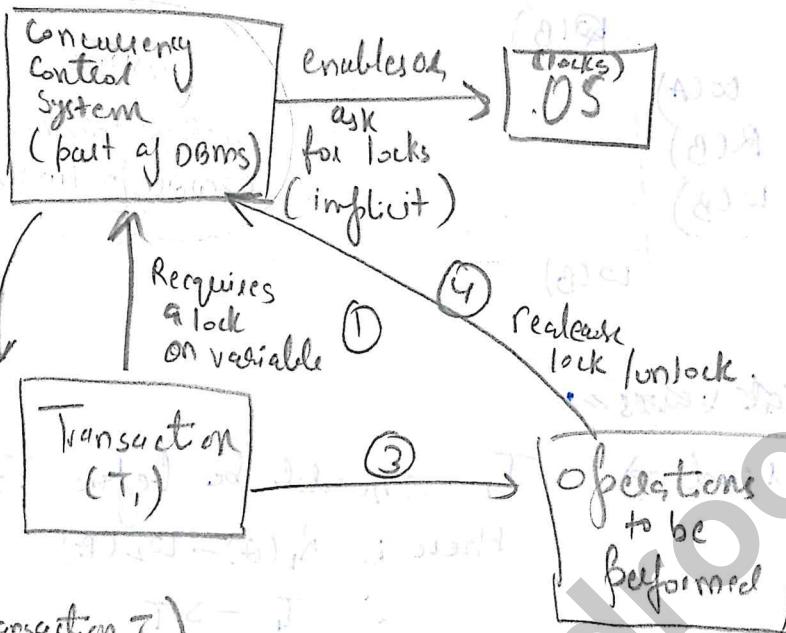


Fig: How locks work.

locks → Shared lock (read lock)
(Only read).

locks → Exclusive lock (write lock)
(Read + write)

Shared lock

Exclusive lock

* A transaction T_1 having shared lock can only read the data.

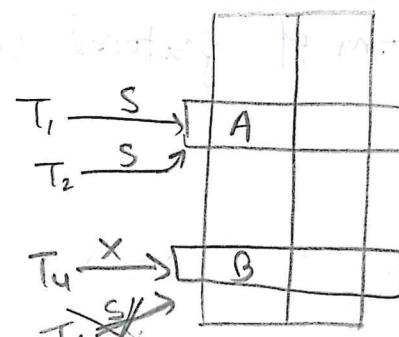
* More than one transaction can acquire shared lock on same variable.

* A transaction T_2 having exclusive lock can both read as well as write.

* Only one transaction can acquire exclusive lock on a variable at a time.

Note - A transaction T_1 , having exclusive lock on any variable will not allow any other Transaction T_2 to even have a shared lock on some variable. It will be denied.

T_3 is denied because it will be denied because T_1 and T_2 are having shared lock on A (\$)

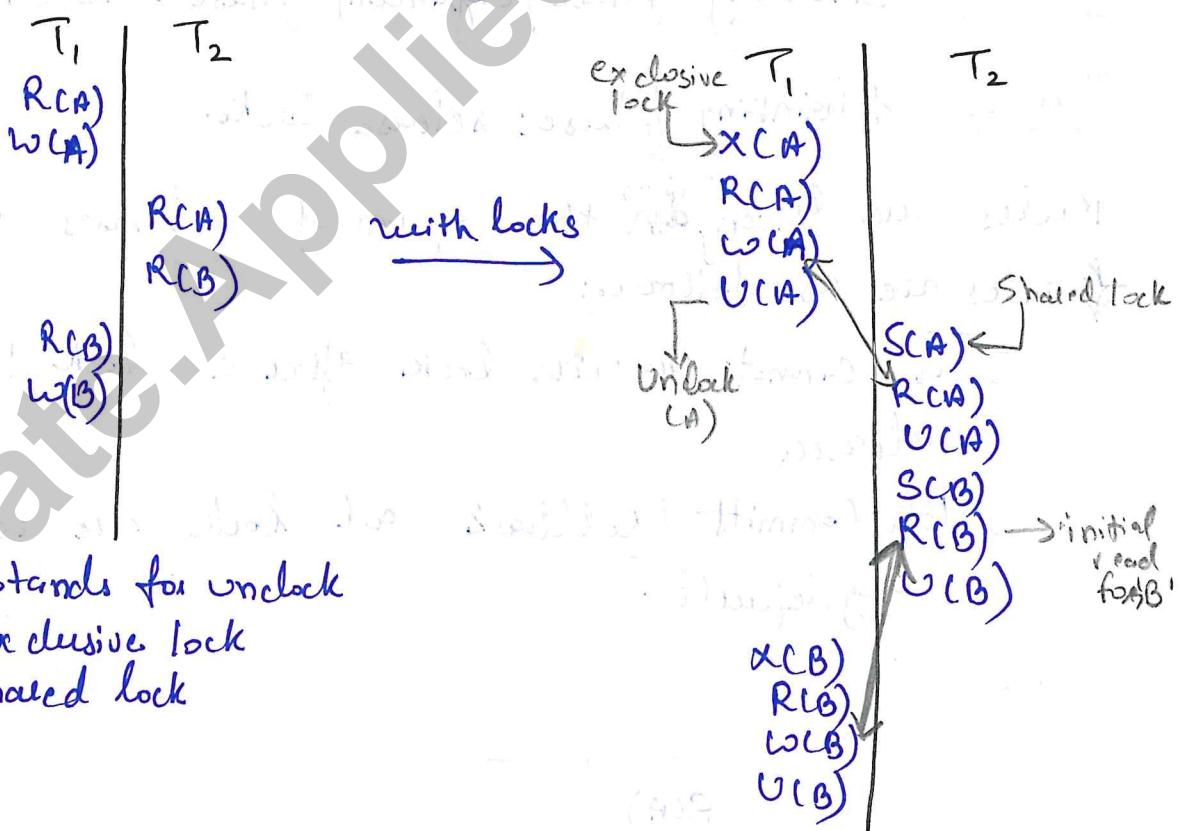


denied lock
and not T4 is having lock (X)
exclusive
on B

So one way to achieve consistency and serializability is through locks.

Example:-

Ex:-



$U()$ → stands for unlock

$X()$ → exclusive lock

$S()$ → shared lock

Here, neither $T_2 \rightarrow T_1$ nor $T_1 \rightarrow T_2$ is valid. So,

locks do not guarantee serializability

The goal/aim behind using locks is to get both consistency and serializability.

In the last example we neither achieved consistency nor serializability.

Therefore, in order to overcome this a set of certain rules are given in the form of protocol called 2 phase locking protocol.

→ 2 phase locking protocol :- (2PL)

It guarantees conflict serializability in turn guarantees serializability.

Conflict serializability \Rightarrow Serializability

This locking protocol has 2 phases:

Phase 1 :- Growing phase / expanding phase : acquires lock

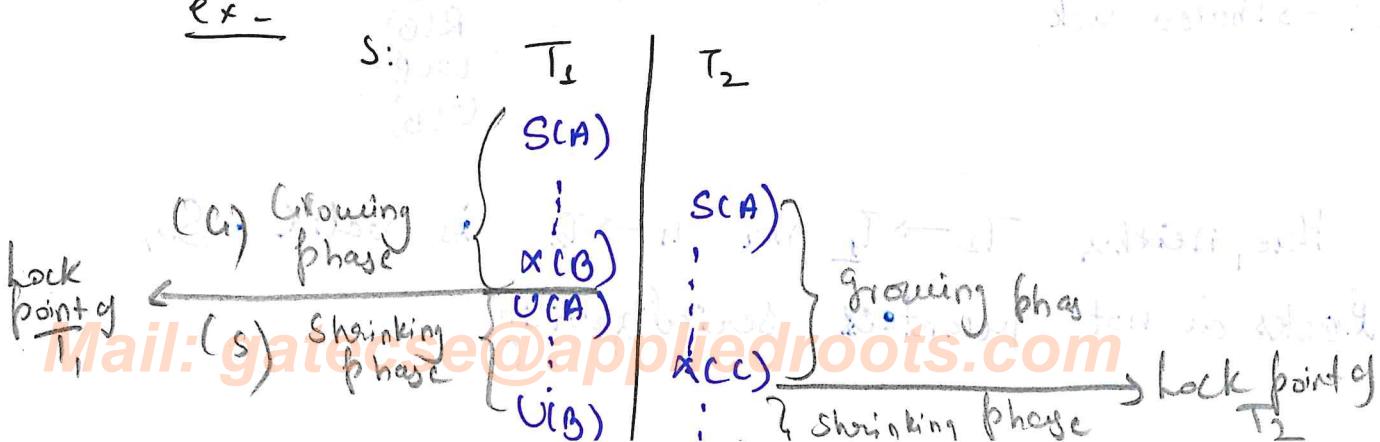
Phase 2 :- Shrinking phase : release lock.

Rules are given in the form of 2 phases where rules are as follows:

[1] T_i cannot acquire lock after a lock has been released

[2] On Commit / rollback all locks are released by default.

Ex -



NOTE - lock point is the point in a transaction where transaction has successfully acquired all the locks and has started releasing the locks.

NOTE - 2PL guarantees CSS \Rightarrow Serializable

Proof:- It is proved based on induction on # transactions

Again CSS $\not\Rightarrow$ 2PL (every CSS do not follow 2PL)

(Q) How to obtain the equivalent serial schedule?

\Rightarrow lock points are used to obtain the equivalent serial schedule. In previous example, lock point of T₁ is followed by T₂ \therefore its equivalent schedule (serial) is T₁ \rightarrow T₂

Example

T ₁
X(A)
R(A)
W(A)
X(B)
V(A)

SC(A)
R(A)
SC(B)
R(B)

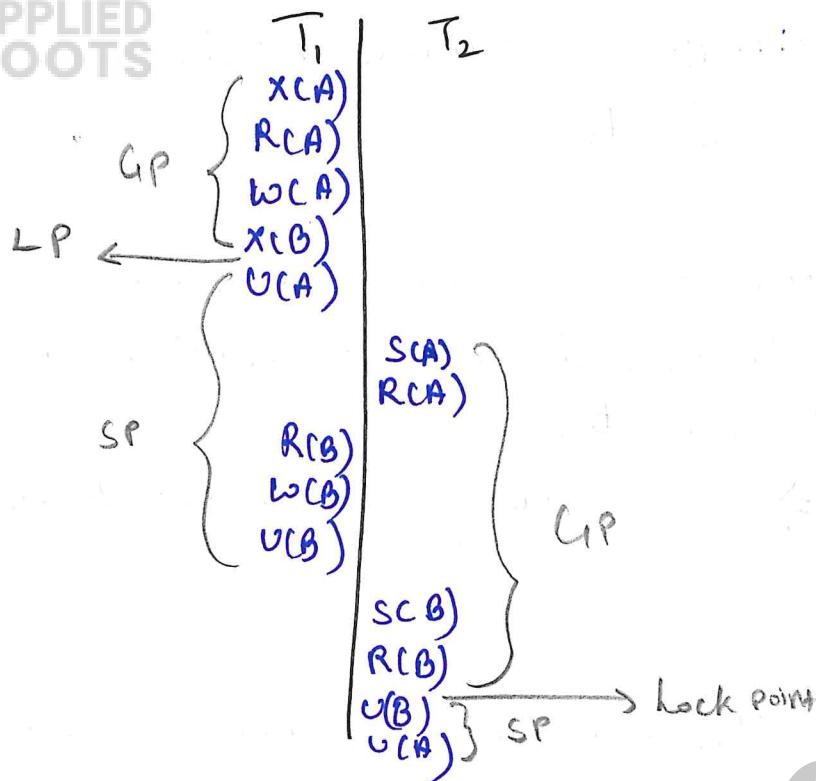
R(B)
W(B)
V(B)

\leftarrow This will be denied bcoz T₁ already have exclusive lock on B.

(Q) Is it 2PL?

Phone: +91 844-844-0102

APPLIED
ROOTS



→ To check for 2PL, we need to discover growing phase, shrinking phase and Lock points.

As we can see, it is following properties of 2PL.

∴ it is CSS and it is equivalent to serial schedule $\Rightarrow T_1 \rightarrow T_2$ (based on the occurrence of these LP)

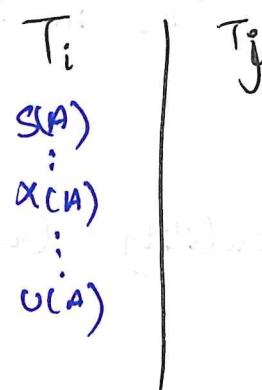
→ Lock upgrade:-

We can upgrade the shared lock to exclusive lock without unlocking it before upgrading.

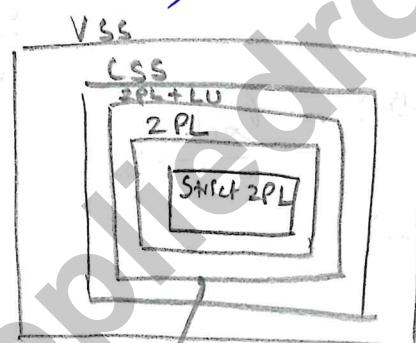
Shared lock $\xrightarrow{\text{upgrade}}$ Exclusive lock

[This is possible only when T_2 is the only transaction having shared lock on it]

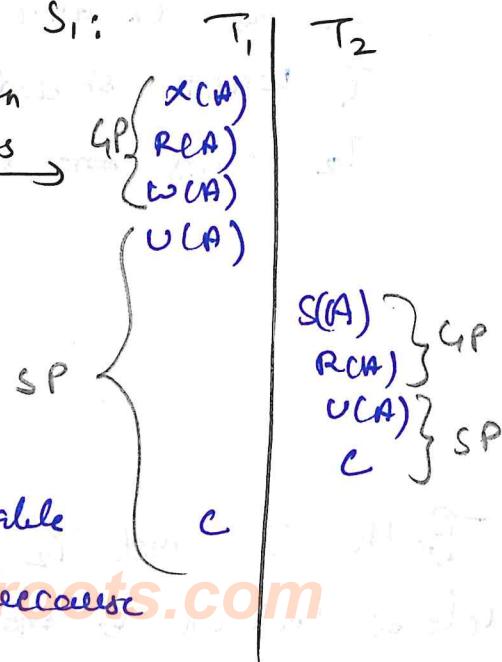
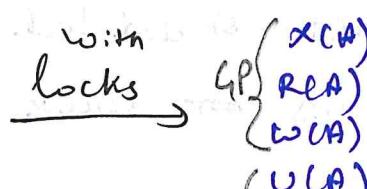
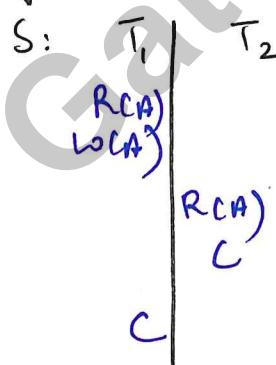
Mail: gatedesk@appliedroots.com

Ex -

Initially T_i has shared lock on A, we can upgrade it to exclusive lock as no other transaction T_j have shared lock on it. (while upgrading there is no need to unlock the shared lock)

NOTE -

2PL + Lock upgrade

Example:-

here, S₁ is 2PL \Rightarrow LSS \Rightarrow serializable

But there is no recoverability because T₂ committed before T₁.

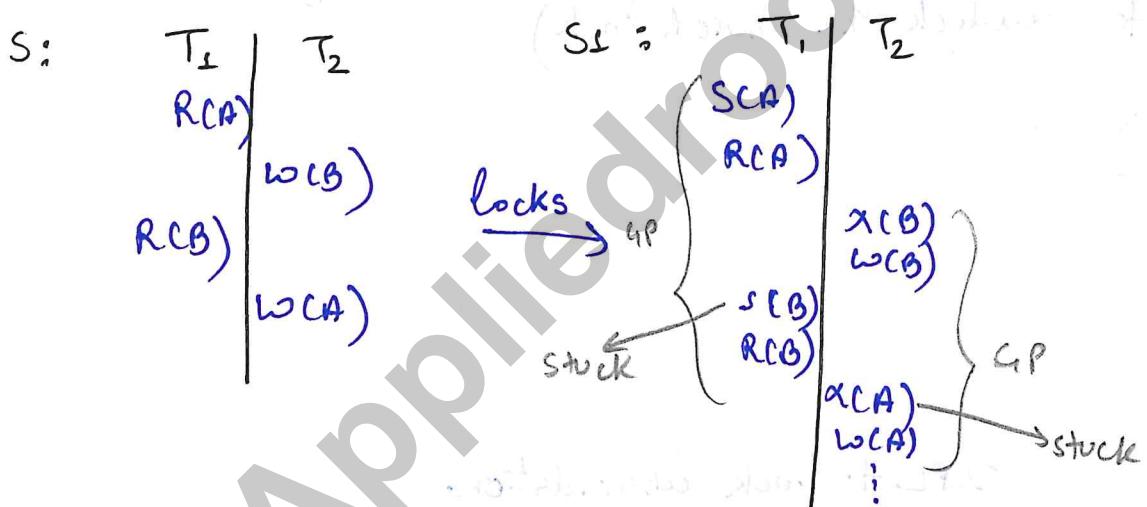
We can't rollback transaction T_1 in $S_1 \because S_1$
is not recoverable

2PL doesn't guarantees recoverability but it guarantees serializability.

Hence, 2PL \Rightarrow Recoverability

→ Deadlocks :-

Q:-

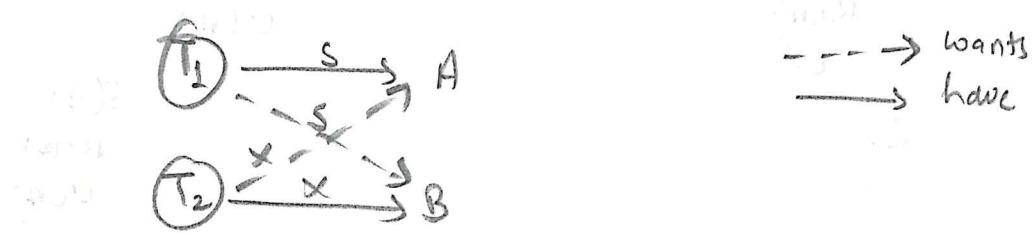


here T_1 has shared lock on A

T_2 has exclusive lock on B

T_1 wants shared lock on B

T_2 wants ~~shared~~ exclusive lock on A



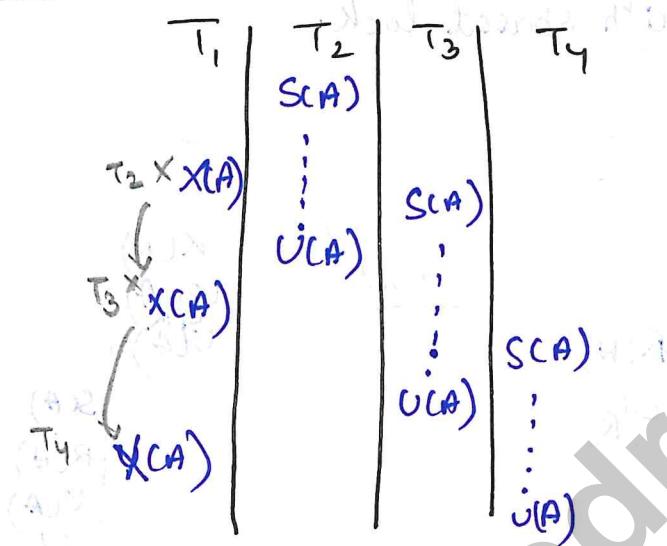
Both T_1 and T_2 are waiting for each other

Mail to gatocsce@appliedroots.com

Neither T_1 nor T_2 can proceed anywhere. This is called deadlock.

This schedule S_1 follows 2PL but it is having deadlock so, 2PL can suffer with deadlocks. (Stuck in infinite loop)

→ Starvation:



Here T_1 is being starved for getting lock on A because T_2, T_3, T_4 are requesting for shared lock (S(A)) before unlocking. This problem is called starvation problem as T_1 is being starved by other transaction. Though it is following 2PL but still it faces the starvation problem.

Note - 2PL suffers from deadlock and starvation.

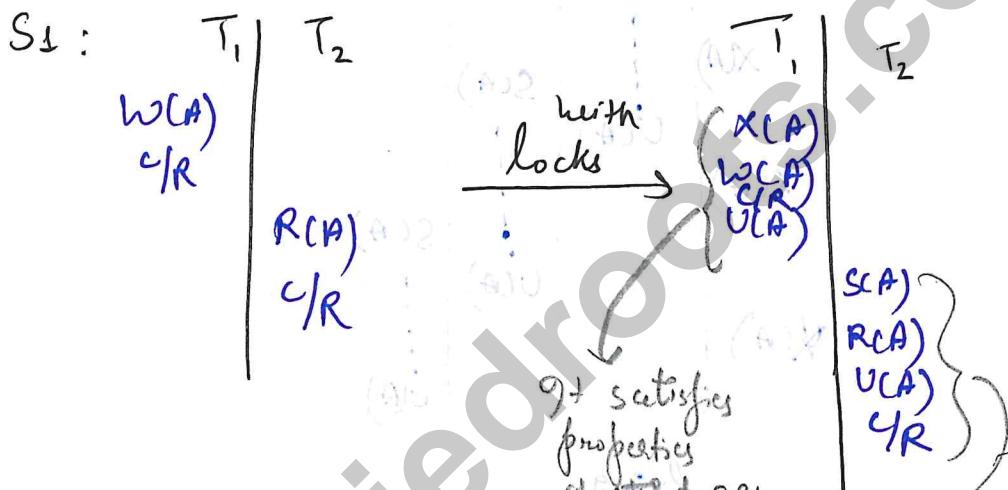
Q) Can we have system which can guarantee both serializability and recoverability?

⇒ Strict 2PL guarantees CS and strict recoverable schedules.

Strict 2PL guarantees both Serializability (SS) and Recoverability (strict recoverable).

If follows all rules of 2PL + all exclusive locks must be held till Commit/rollback but there is no such restrictions with shared lock.

Ex -



It satisfies properties of strict 2PL

After the commit

was not there

than it would

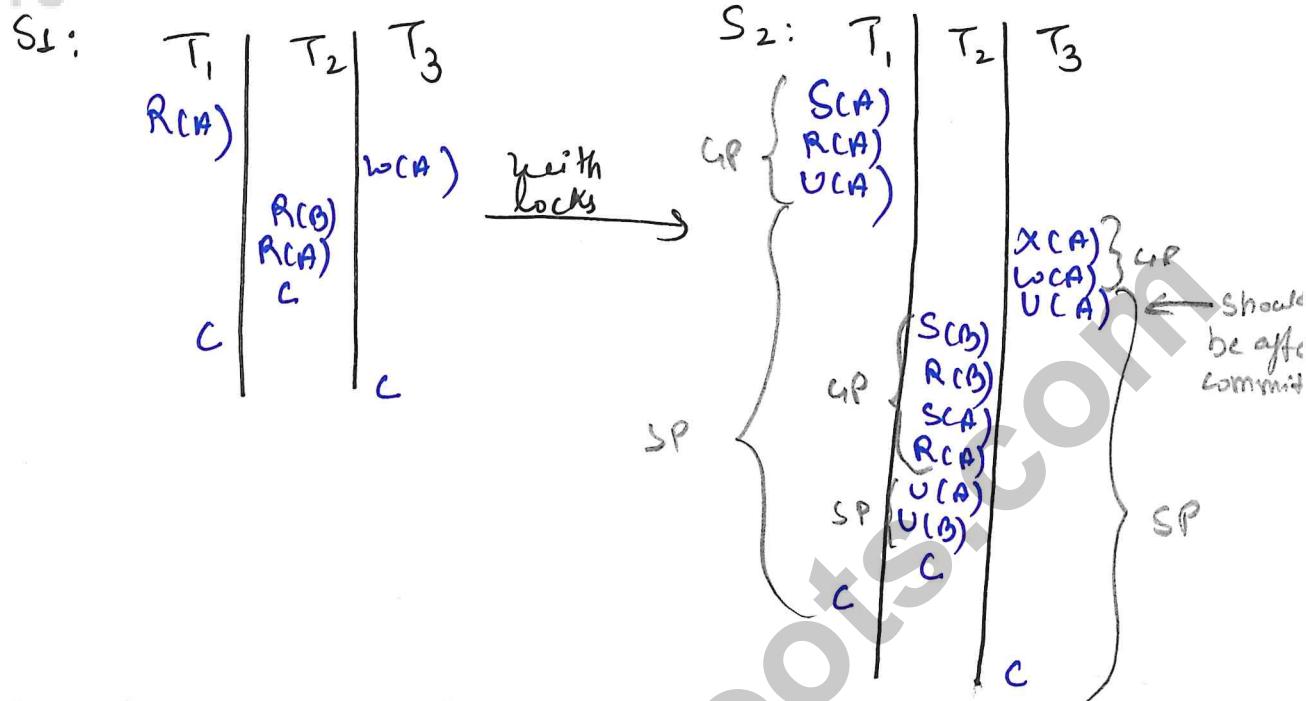
not be strict 2PL

Strict 2PL \Rightarrow (SS + strict recoverable)

(Q) How to obtain equivalent serial schedule given that schedule is having Strict 2PL?

→ The sequence of lockpoints will be the order of transaction. The resultant order will be the equivalent serial schedule.

(Q1) Is it strict 2PL?

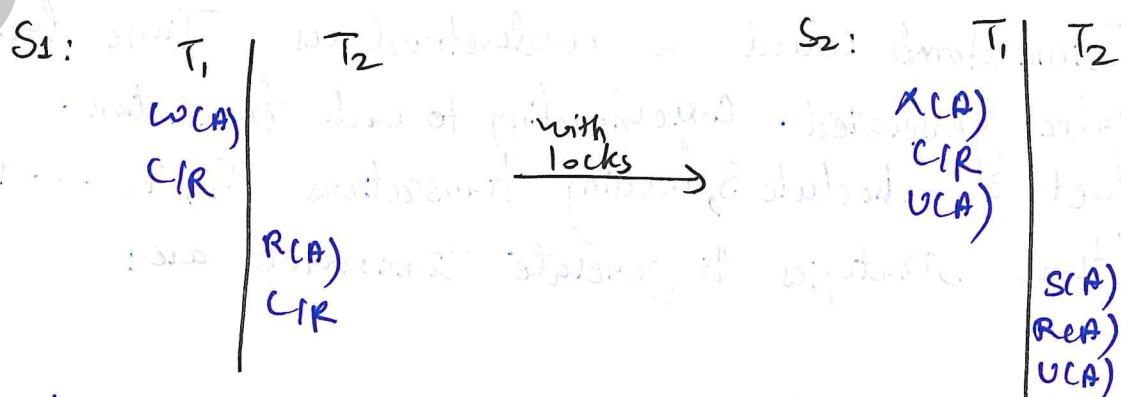


→ S_2 follows all properties of 2PL but it is not strict 2PL because in T_3 before committing after $w_3(A)$ we are releasing lock.
 $\therefore S_2$ is not in strict 2PL.

→ Strong 2PL / Rigorous 2PL

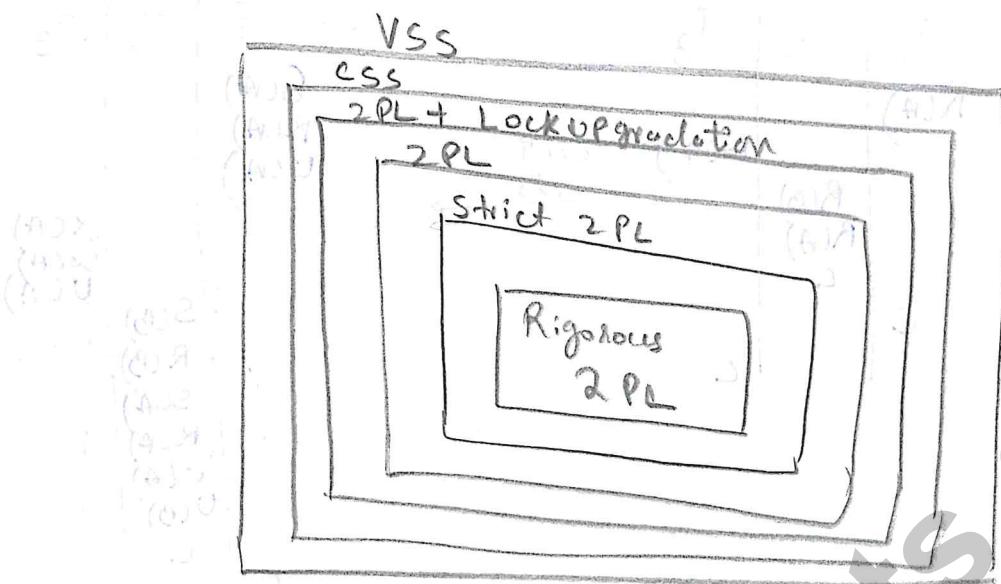
Strong 2PL + all locks (shared or exclusive) should be held till transaction commits or rollback.

Ex -



S_2 is having exclusive lock and it is released after Commit. S_2 is Strict 2P. But it is not Rigorous 2PL because it releases lock before commit.

NOTE - Strict 2PL may suffer with deadlock and starvation.



18.0 Timestamp based protocols + Deadlock and starvation prevention:

Another way to strategy to obtain both serializability and recoverability is Time stamp based protocols.

In lock based protocol system we need to have OS enabled locks / software locks while Timestamp based protocols are alternative of lock based protocols.

Time stamp could be understood as Time based value generated corresponding to each transaction.

Let a schedule S, having transactions $T_1, T_2 \dots T_n$, The two strategies to generate timestamp are:

Generate a time stamp:

- ① Use the system time generated by OS to generate TS (time stamp)
- ② Create a shared counter and use it to assign a unique no. for each transaction.



This counter is thread safe.

→ Time stamp based protocol: (TSP, TSOP)

Timestamp ordering protocol

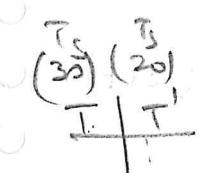
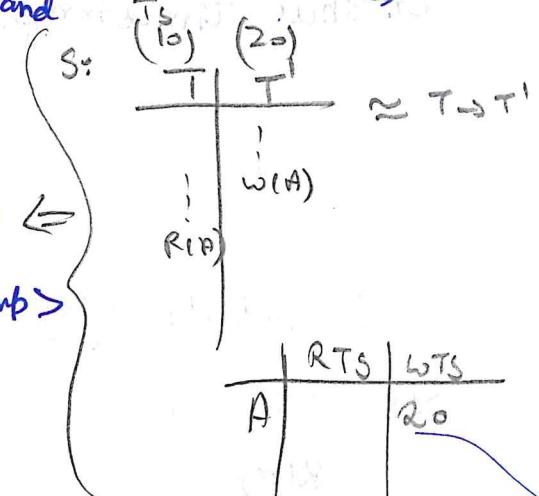
- * Each transaction must obtain a unique timestamp @ start.
- * Each data-item also has time stamp (Write time stamp (wTS), Read Time stamp (RTS))

Given the above construct the following 2 conditions should be satisfied.

Case 1: Let assume transaction T and T' (T locates to R(A))

$$[1] \quad TS(T) < wTS(A)$$

→ Rollback (T) and restart with new timestamp (new timestamp > T)



$$[2] \quad TS(T) > wTS(A)$$

→ execute R(A)

$$\rightarrow RTS(A) = \max(RTS(A), TS(T))$$

RTS	wTS
30	20

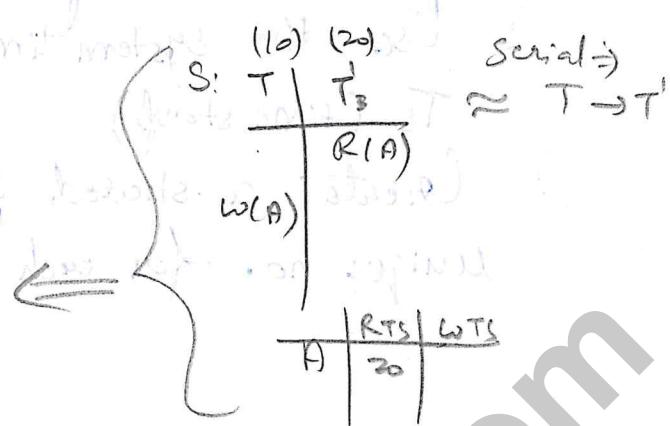
Case 2:



Want to $w(A)$:-

(1) $RTS(A) > TS(T)$

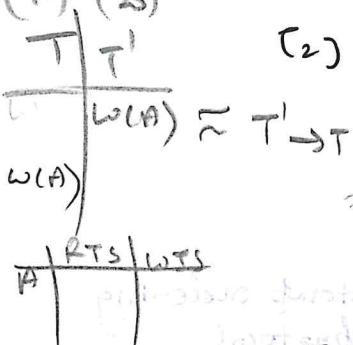
→ Roll back (T) and restart



(1) (2)

(2) $WTS(A) > TS(T)$

→ Roll back (T) and restart



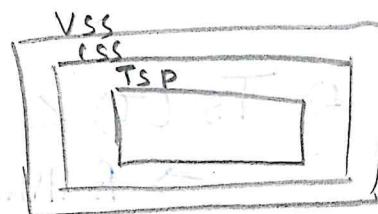
[3] otherwise,

→ execute $w(A)$

→ set $WTS(A) = TS(T)$

NOTE:-

* Time stamp based protocol guarantees CS (no clity reads). Given any TSP, we can obtain equivalent serial schedule by simply ordering transactions based on their timestamps.



→ Deadlocks w.r.t TSP :-

S1:	T_1	T_2
	$R(A)$	
		$w(B)$

S1 suffers with deadlock

in case of 2PL, strict 2PL.

Let's check it with respect to (wrt)

TSP.

Phone: +91 844-844-0102

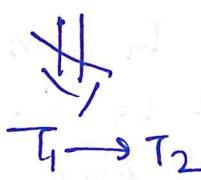
Since T_L has smaller T_S , \therefore it should be equivalent to serial schedule

$$T_1 \rightarrow T_2$$

Since S_L is not equivalent to

$T_1 \rightarrow T_2$ as TSP will not allow $R_1(B)$ to take place.

It will say rollback ~~on~~ $R_1(B)$ and restart the whole thing.



Hence, * TSP guarantees no-deadlock. It does not allow deadlocks to take place by not executing the operation leading to deadlock.

→ Recoverability w.r.t TSP:

	TS: 10	20
S_L :	T_1	T_2
	$w(A)$	$RCA)$
	$C(R)$	C
	RTS	WTS
A	(20)	10

it is an irrecoverable schedule

TSP will try to check whether Schedule S_L is equivalent to serial schedule $T_1 \rightarrow T_2$ (as T_1 has smaller timestamp).

TSP will execute $R_2(A)$ ~~on~~ as $TS(T) > WTS(A)$ $\therefore R_2(A)$

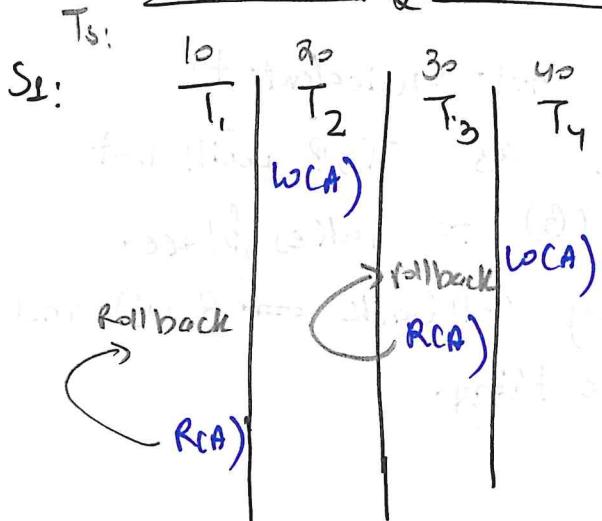
will get executed and RTS will be updated to 20.

\therefore TSP is allowing a irrecoverable schedule



Hence, TSP does not guarantee recoverability

→ Starvation wrt TSP:



Task of TSP, to make S₂ equivalent to $\approx T_1 \rightarrow T_2 \rightarrow T_3 \rightarrow T_4$

Here T₂ and T₄ com starve T₁ and T₃. ∴ starvation com happen as Transactions with larger TS com force Transaction with lower TS to rollback.

∴ TSP doesn't guarantee starvation free schedule

→ Strict Timestamp ordering protocol (Strict TSP)

TSP ensures serializability and no-deadlocks while strict TSP ensures serializability + no-deadlocks + strict recoverability. Basic construct is same but in place of 2 cases (as in TSP), strict TSP has only 1 case. (Because for strict recoverability we can't have WR and CW)

So, if any transaction T wants to read or write
ie $R(A)$ or $W(A)$

if, $WTS(A) < TS(T)$

\rightarrow Rollback (T) and restart

Rule for
strict
TSP

else,

Everything stays the same, just update the
WTS and RTS accordingly

This ensures No WR and No WW and hence, ensures
strict recoverability. Since, this condition is more
strict than TSP \therefore it definitely ensures no-deadlock and
Serializability.

NOTE

Strict TSP \subset TSP \subset CSS \subset VSS

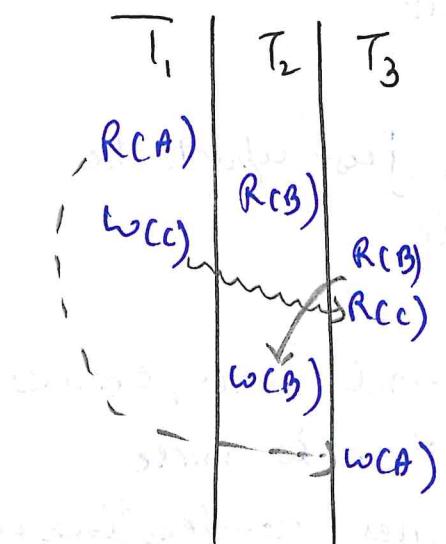
(it is also
strict
recoverable)
No-deadlock
but starvation
may occur.

(Q1) What are the valid $TS(T_i)$ values to avoid
rollbacks?

SL:

T_1	T_2	T_3
$R(A)$		
	$R(B)$	
$W(e)$		
	$W(B)$	
		$R(e)$
		$W(A)$

\Rightarrow In order to assign TS we need to find the given schedule S is equivalent to which serial schedule.



T_1 should happen before T_3 due to $R_1(A)$ and $W_1(C)$. $T_1 \rightarrow T_3$

T_3 should happen before T_2 due to $R_3(B)$. $T_3 \rightarrow T_2$

i.e. Serial order is

$$T_1 \rightarrow T_3 \rightarrow T_2$$

So if we assign TS values as

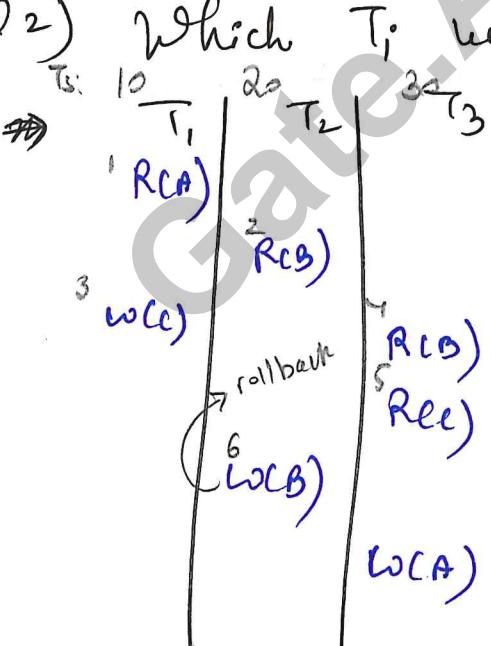
$$T_1 = 10$$

$$T_3 = 20$$

$$T_2 = 30$$

These, TS values will ensure that if we apply TSP then there will be no rollbacks.

(Q2) Which T_i will rollback first? (TS are already given in question)



\Rightarrow For the give TS values

	RTS	WTS
A	10	
B	20	30
C	30	10

i.e. the equivalent serial order TSP trying to maintain's $T_1 \rightarrow T_2 \rightarrow T_3$ (acc. to TS)

while performing 6th operation

$$TS(T) = 20, RTS(B) \leq 30, WTS(B) = 20$$

and equivalent schedule is $T_1 \rightarrow T_2 \rightarrow T_3$

$$\text{so } \text{aff}^o(T_3) < TS(T_3) \quad (RTS(B) \geq WTS(B))$$

(a) since $W_2(B)$ will not be allowed and it will be rolled back.

\therefore First transaction which will roll back is T_2 .

NOTE- In real world combination of 2PL + TS are used, as

2PL } Cannot avoid deadlock + starvation
Strict 2PL }

TS } Can avoid deadlock but not starvation.
Strict TS }

In real world combination of 2PL and TS are used, as TS can effectively avoids deadlocks.

\therefore Strict 2PL + TS are used in practice to ensure serializability + recoverability + no-deadlock.

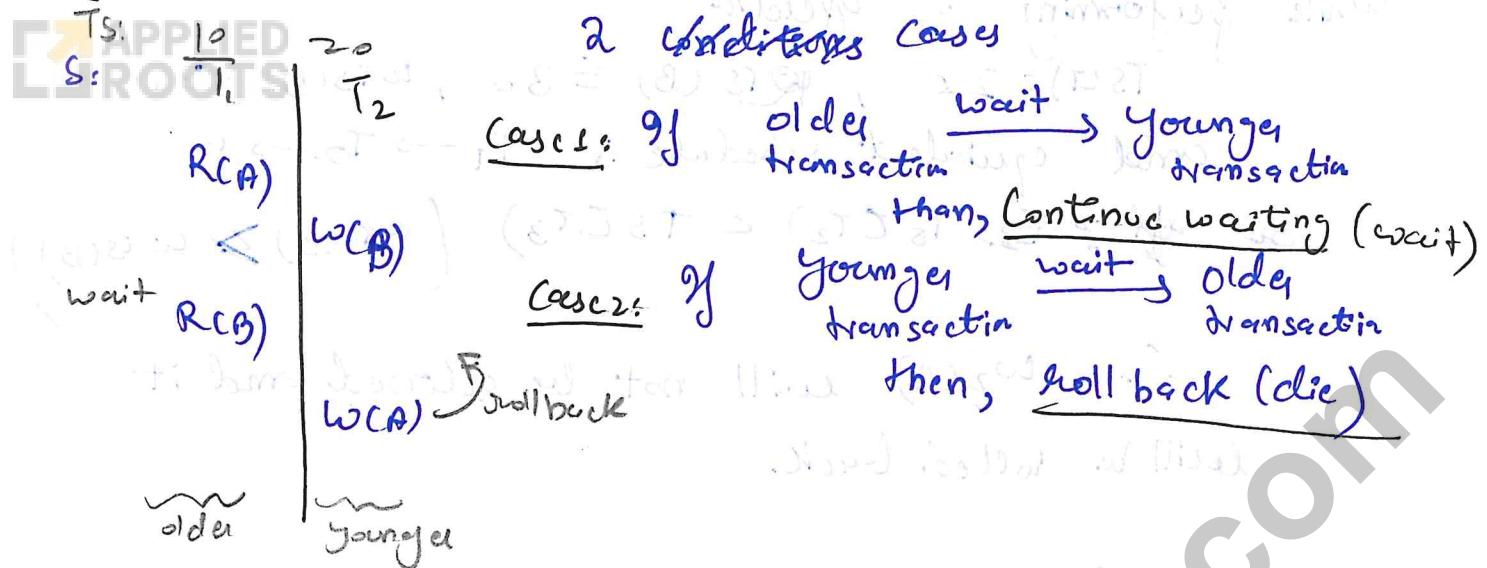
(Q) Additional strategies to be introduced in lock based protocol to ensure deadlock and starvation?

\Rightarrow There are 2 strategies

(1) wait-die

(2) wound-wait

(1) Wait-die : Use strict 2PL + TS with



short Note that T_1 is holding Shared Lock on A while T_2 have exclusive lock on B.

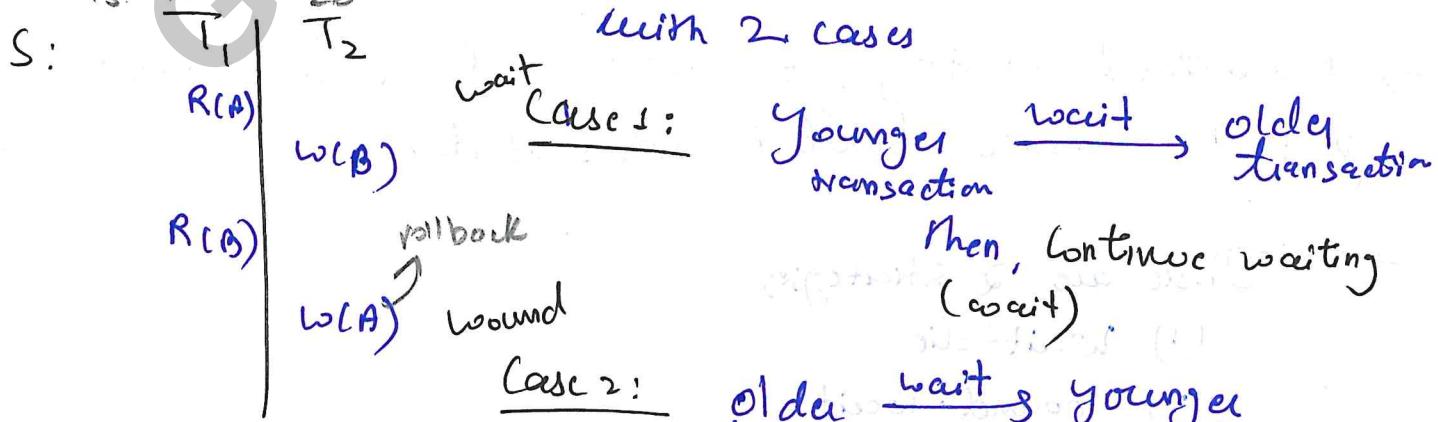
$\therefore T_1$ is older transaction waiting for T_2 (younger transaction) to release lock on B. \therefore acc. to rules, it will keep waiting.

T_2 is the younger transaction which is waiting for T_1 (older transaction) to release lock on A

\therefore Based on rules, it will get rolled back (Die).

Hence, S2PL+TS will overcome the deadlock problem.

(2) Wound - wait : It also uses strict 2PL + TS



Mail: gatecse@appliedroots.com

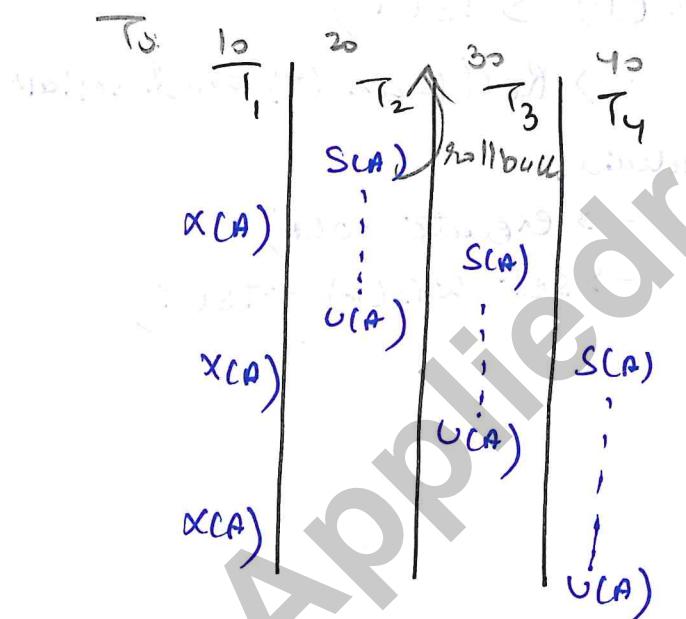
NOTE

In both the cases wait-die and wound-wait we rollback the younger transaction.

→ Starvation prevention:

* Wound-wait:

- ① Older $\xrightarrow{\text{wait}}$ younger $\xrightarrow{\text{rollback}}$
- ② Younger $\xrightarrow{\text{wait}}$ older (Waiting)



With the help of wound-wait we can prevent starvation.

Hence, Strict 2PL + wound-wait (TS) will ensure

Strict 2PL guarantees this	\rightarrow Serializability	} Strict TSP guarantees this
	\rightarrow Strict Consistency	
Wound-wait(TS) ensures this	\rightarrow no-deadlock	}
	\rightarrow no-starvation.	

It is based on Time stamp based protocol.
A simple TSP have 2 cases. The 2nd case of TSP has 3 conditions

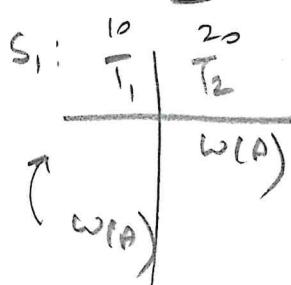
Case 2: — T wants to write (A) / w(A)

$$[1] \quad RTS(A) > TS(A)$$

→ Rollback (T) and restart

$$[2] \quad WTS(A) > TS(T)$$

→ Rollback (T) and restart



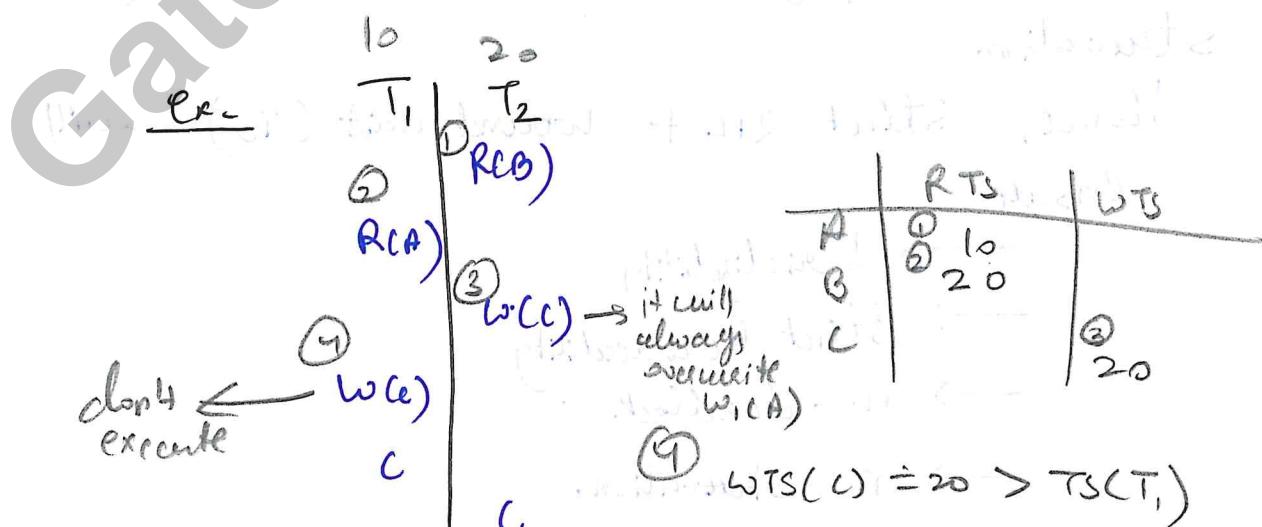
[3] Otherwise,

→ execute w(A)

→ set $WTS(A) = TS(T)$

This rule is changed
in Thomas
write rule

It often refuses to
as Thomas
write rule.



Phone: +91 844-844-0102

According to serial schedule, anyway T_2 's write will overwrite the T_1 's write. (anyhow $w_i(A)$ will be overwritten). So TWR say if Transaction T_2 has larger TS and has written a value of variable and we have an older transaction T_1 , who is writing after it (in serial order anyway T_2 will overwrite T_1), then don't execute $w_i(A)$.

So, it says don't rollback # if such condition occurs, just ignore that operation.

TWR is used in context of TSP, it is not used in isolation.

↳ In general, standard SQL database allows a local transaction to commit and affect a part of database while other parts of the same transaction remain uncommitted. This is called two-phase commit protocol. It consists of two phases:

- Preparation phase: All changes made by the transaction are recorded in log files.
- Commit phase: If no errors occur, then all changes recorded in log files are made permanent.

APPLIED STRUCTURES (Sequential files, Phone: +91 844-844-0102)

indexing, B and B⁺ trees)

19.1 files and indexing: Introduction:

DB-tables storing Hard Disk (Secondary memory)

Objective: Access data on the disk fast
→ search, read, write, modify, delete

→ why do we need index?

Index makes searching of an attribute / column or a tuple more faster.

Index is a data structure such that access based on that attribute is faster. (It is suggested to have indices on primary key).

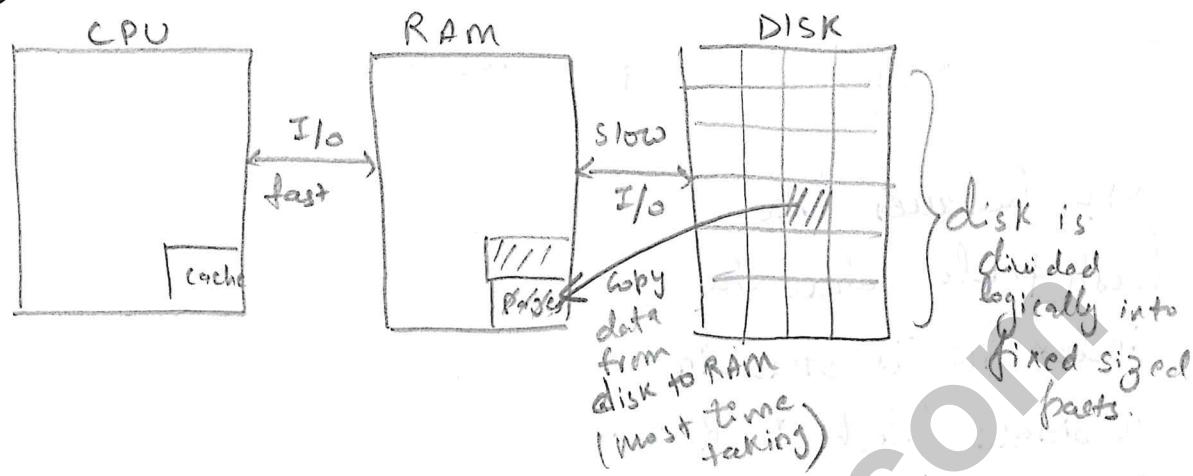
In the given example id is the index + search by name.

<u>id</u>	name	addr
num	str	str
1	n ₁	a ₁
2	n ₂	a ₂
3	n ₃	a ₃
:	:	:

(index)

Mail We can have multiple indices in a given table.

Accessing data on disk:

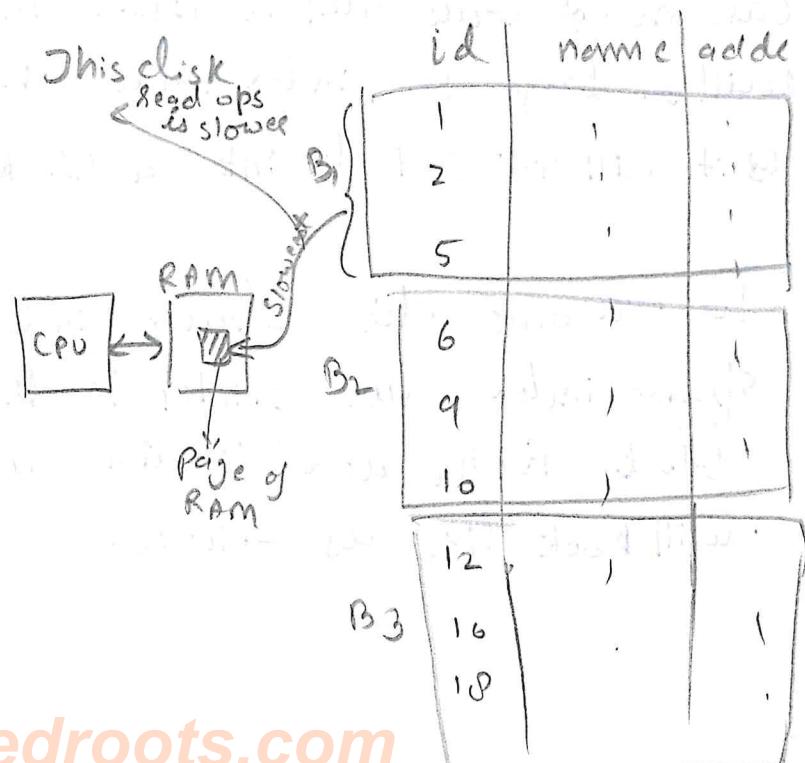


So, how much time it takes us to execute an SQL depends on how efficiently we copy data from disk to RAM and back. (Everything is stored in blocks in a disk, a file/table is also divided into blocks and # of records stored in each block depends on size of each record)

If size of each record is S and Block size is 512B then no. of records it can store is $512/S$.

Cost of accessing data:

The DISK-READ operation is very costly, but once the block is loaded in RAM than it is very cheap to read data and perform required operations.



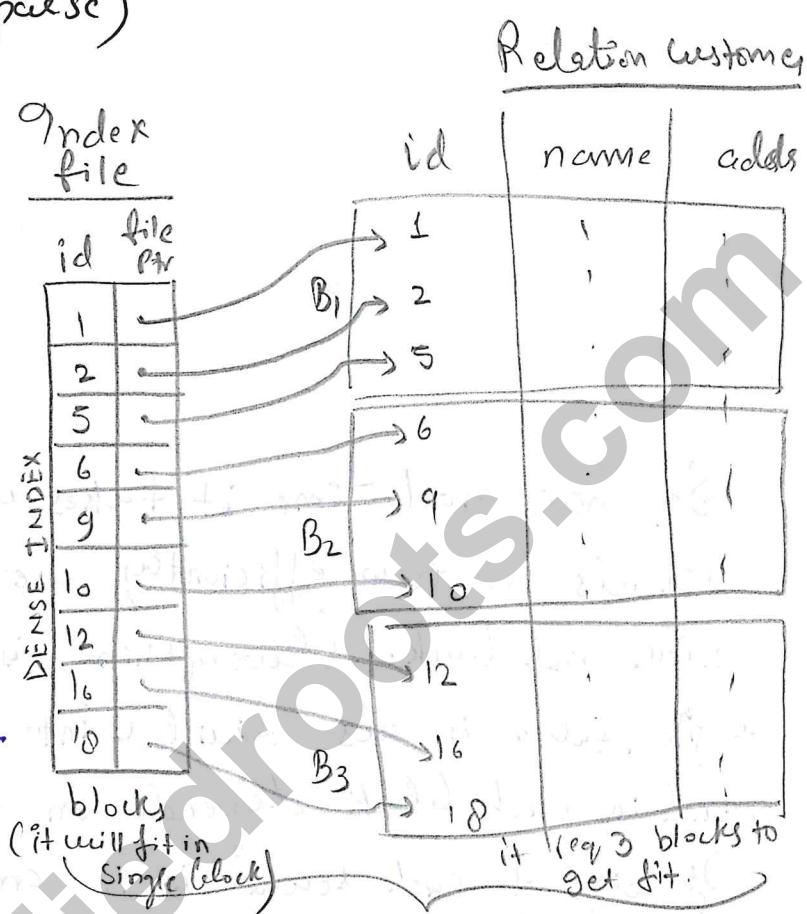
No, our aim is to minimize the number of disk reads.

→ Indexing: (dense) sparse)

It behaves like hash table where the index file is stored in a single block in RAM.

It is dense because for every row we have a pointer.

Sparse indexing is when we have pointers for a subset of rows.



Disadvantage of dense index is when we have large tables. In this case no. of rows will be high. ∴ Index table will be large so, index table will also get divided in blocks as it will not fit into a block.

Both of them are stored in DISK

Fig: Dense indexing

To overcome this sparse index came into picture. Sparse index have pointer to the first row of each block. For the above relation customers : Sparse index will look like as follows.

Our goal is to minimize disk-reads.

The index file will be smaller

in case of sparse index and

\therefore could easily fit in a block

in RAM or disk

So with just 2 disk reads

we can get the required tuple.

Search 10 -
2 disk reads

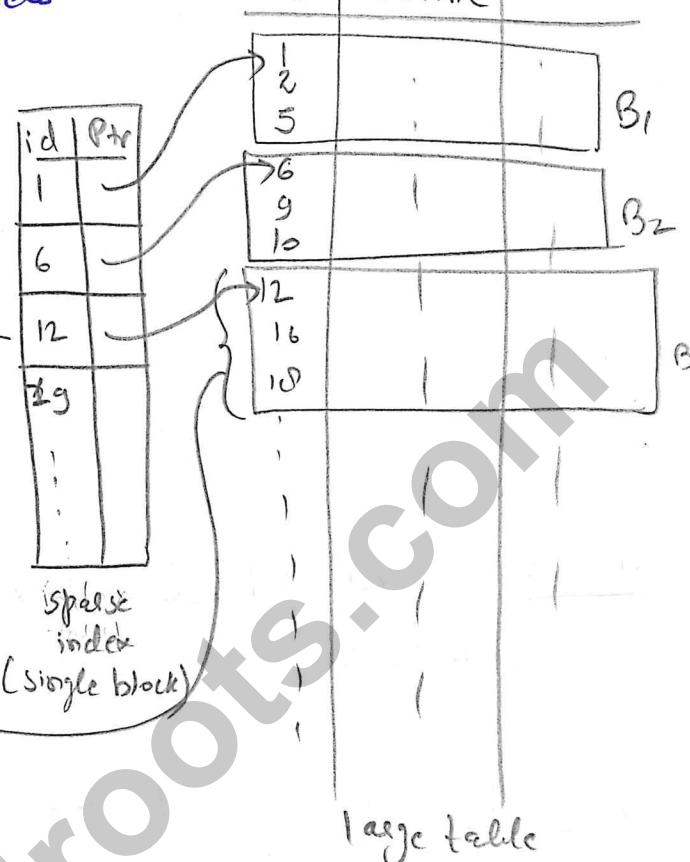


Fig: Sparse indexing

Suppose we have super large table and even sparse index

could not fit in a single block of disk.

In this case we will have.

(m+1) disk reads in worst case

This leads to the idea of
Multilevel indexing

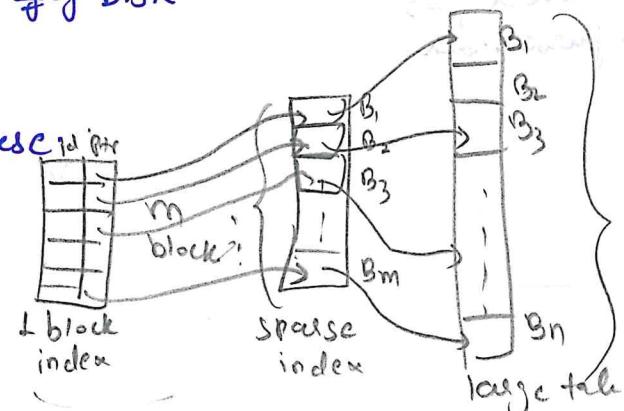


Fig: multilevel index

In multilevel indexing, we build an index out of index so that we can fit the index into a single block of disk or RAM. The multilevel index seems to be like a big tree. We can arrange the index in the form of m-way trees. An m-way search tree will have m keys and m+1 pointers in a node. (Connect to the concept of BST)

m Keys with m+1 Points

Phone: +91 844-844-0102

APPLIED
ROOTS

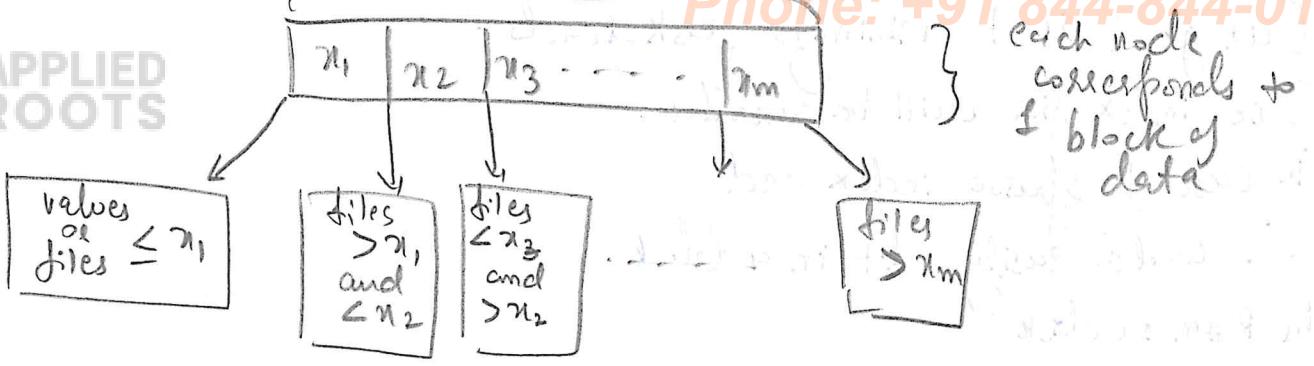


Fig: m-way tree node structure.

NOTE - multilevel index is nothing but m-way search trees.

19.2 B-Trees and B+ Trees with examples:

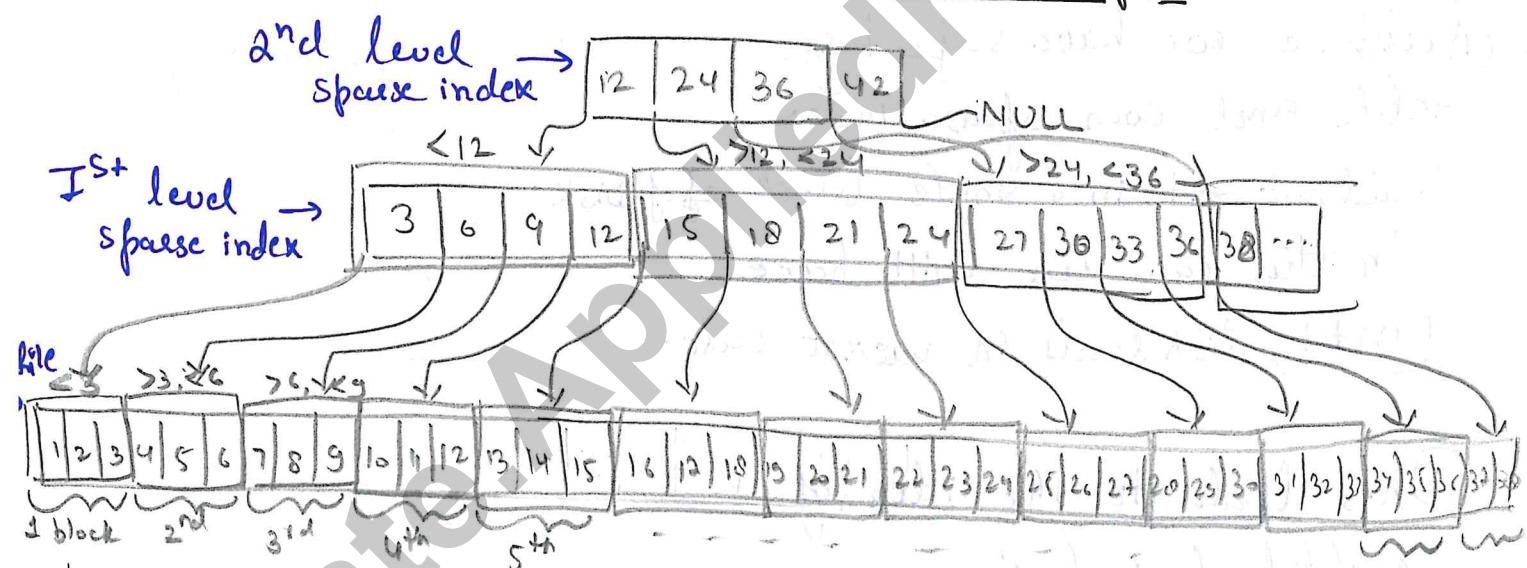


Fig: multilevel index as an m-way search tree

id	name	order
1	n1	a1
2	n1	a2
3	n2	a3

On searching 14, we will have 1st read on 2nd level

2nd read on 1st level index and then final read of record : total 2 reads are required in this case.

→ How to construct an optimal m-way tree?

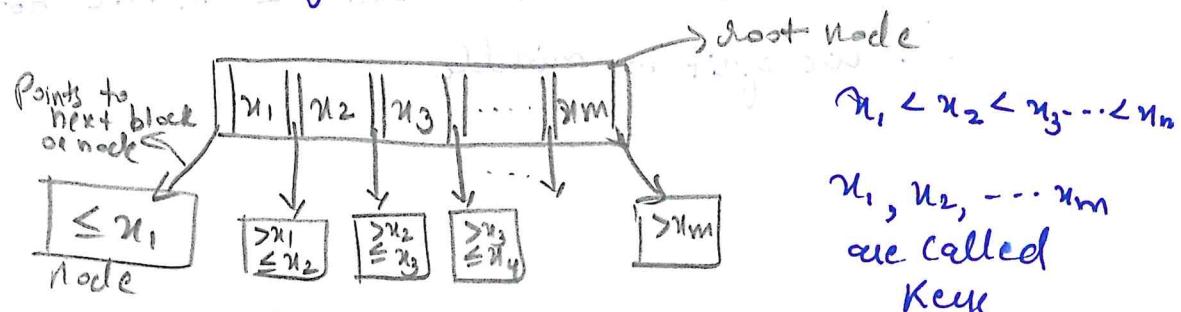
Using m-way trees we have reduced the no. of disk reads. Now the next task is to store them minimum no. of blocks so that height of the tree does not exceed. If each of the node/block is totally filled then we can achieve minimal depth of tree (Pack the data), this is the only way to get optimal m-way tree. (We want height of the tree as small as possible as height of the tree = no. of disk reads in approx).

For this, the concept of B-tree and B⁺ trees come into picture.

* B-trees:

- * Inspired by/m from m-way trees.
- * It is a rooted tree
- * Each node in B-tree is a block

Typical structure of a node is as follows:



Mail: gaurav@appliedroots.com
So each node have m keys and m+1 pointers.

These $u_1, u_2 \dots u_m$ are pointers to a tuple with id $u_1, u_2 \dots u_m$ respectively.

- * All leaf nodes are at some depth.
- * All internal nodes (non-leaf node) except root must have $\lceil m/2 \rceil$ children. (each of nodes must be half filled)
- * Root has atleast 2 children.

These 3 rules ensure that

depth of tree is as small as possible.

- * Order is maximum no. of pointers of a node in a tree.

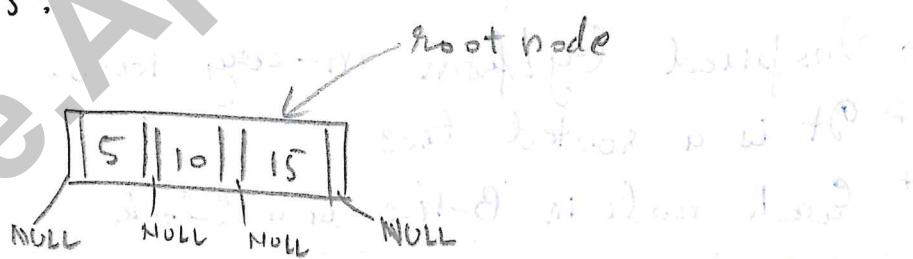
~ Build a B-tree (insert):



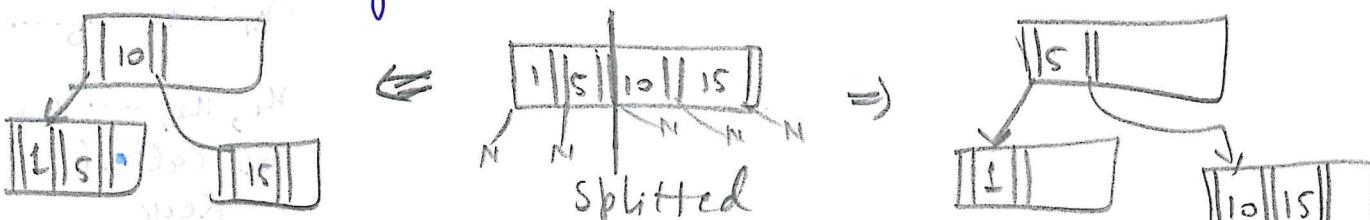
\therefore Order of tree = m+1 = # Ptr in a tree at max

If order = 4, then no. of keys possible in a node = 3

① Insert 5, 10, 15:



② Insert 1: We cannot insert 1 in the root node.
∴ we split in middle

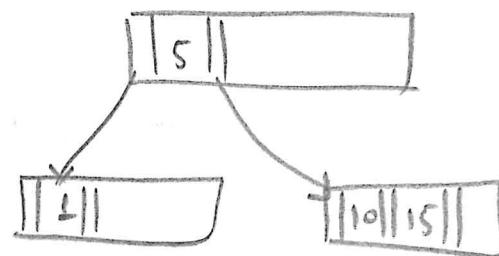


Case 1: Left biasing
(splitting st more elements are in left subtree)

Case 2: Right biasing
(splitting such that more elements are on right tree)

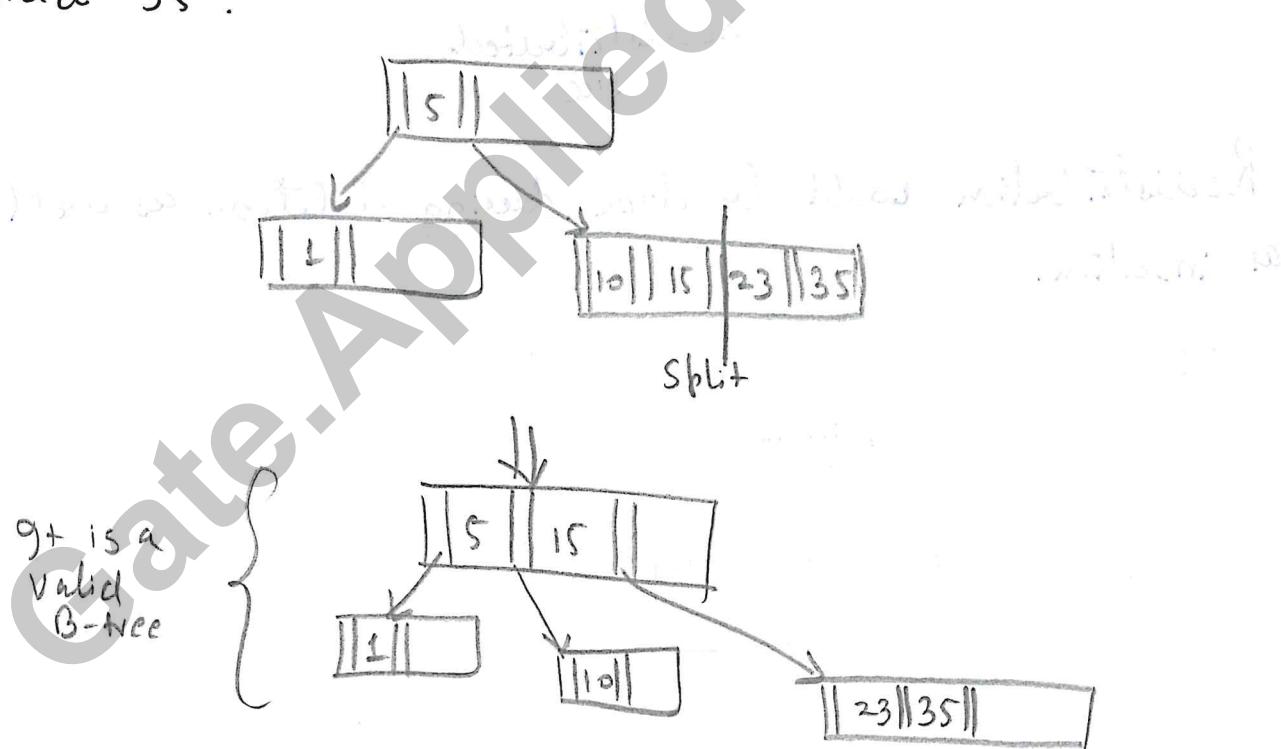
Note that both right biasing and left biasing is valid. Let us continue with right biased tree (ie case 2).

③ Add 23 :



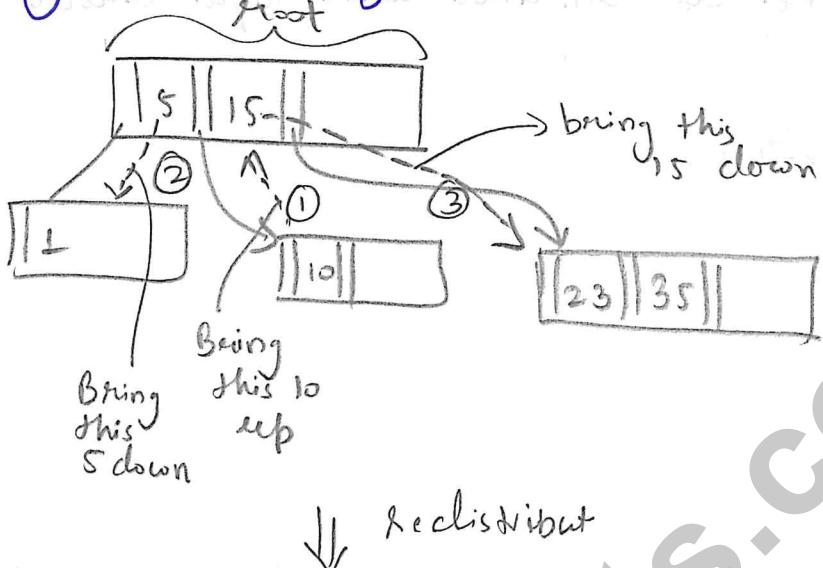
We always fill tree in bottom up method. ∴ we have kept in the node where we have some space.

④ Add 35 :

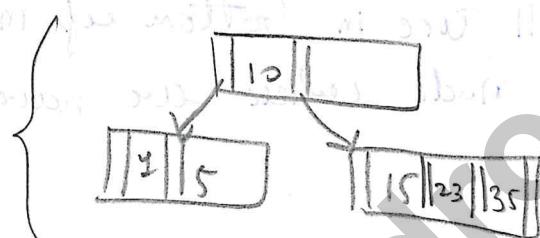


But here, there is some problem as there is lot of empty space in each node. So, we should somehow adjust them so that we can get utilize those empty spaces. ie Reallocate the data/keys.

On redistributing we will get



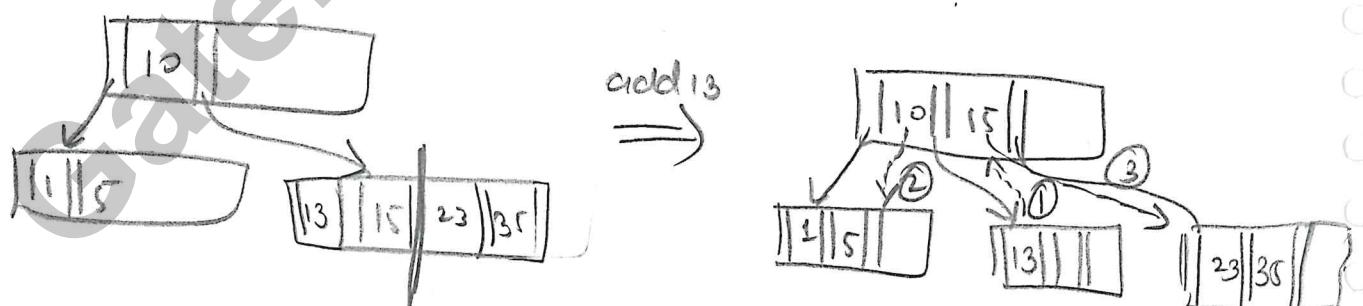
We have saved spaces by making it into 3 blocks in place of 4 blocks.



Redistributed tree

Redistribution could be done during deletion as well as insertion.

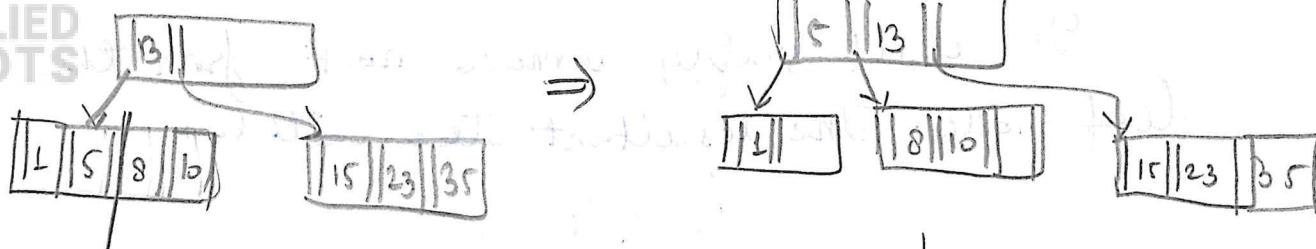
⑤ Add 13:



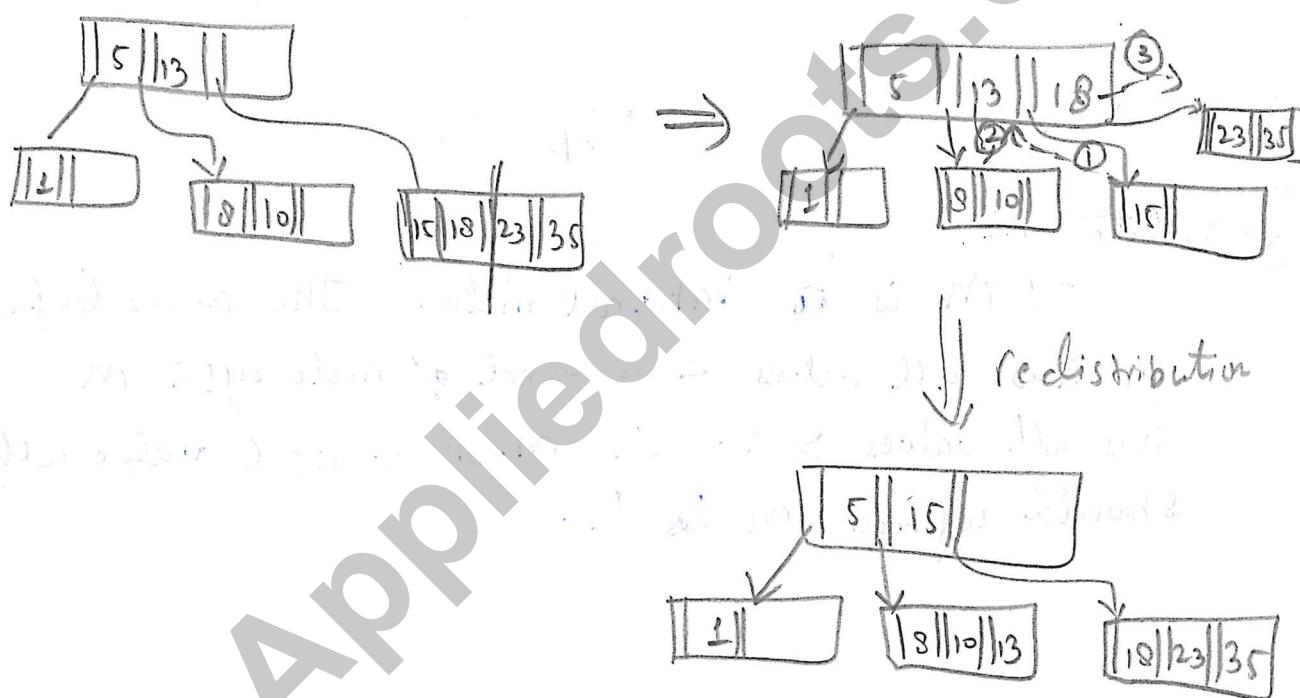
This is not insertion since it will redistribute the process. Insertion can be done if we want to keep the original structure.

⑥ Add 8 :

Phone: +91 844-844-0102 293-



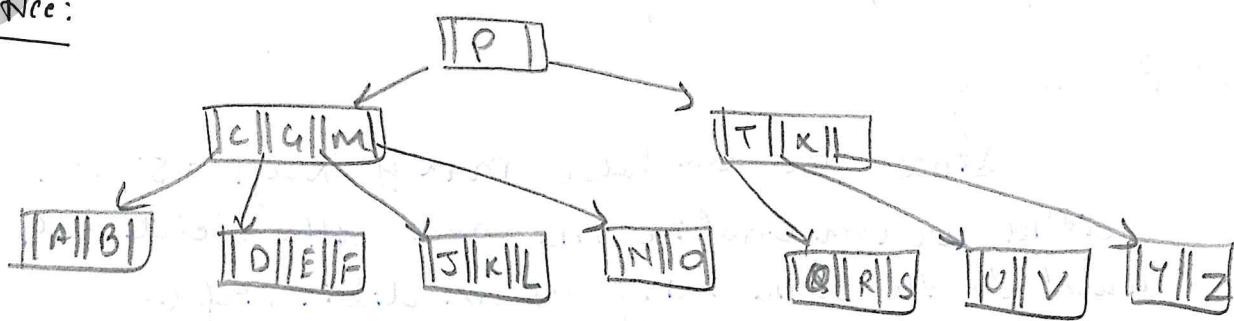
⑦ Add 18 :



~ Deletion from B-tree (Delete) :

Order = 6 \therefore max # Keys = 5

initial tree:

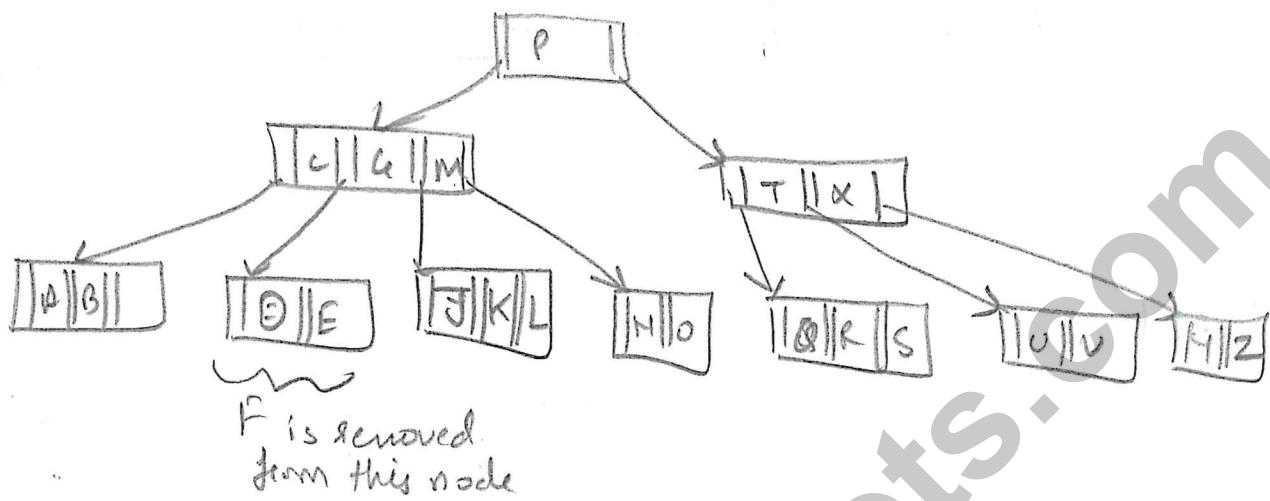


Mail: gatecse@appliedroots.com

① Delete F :

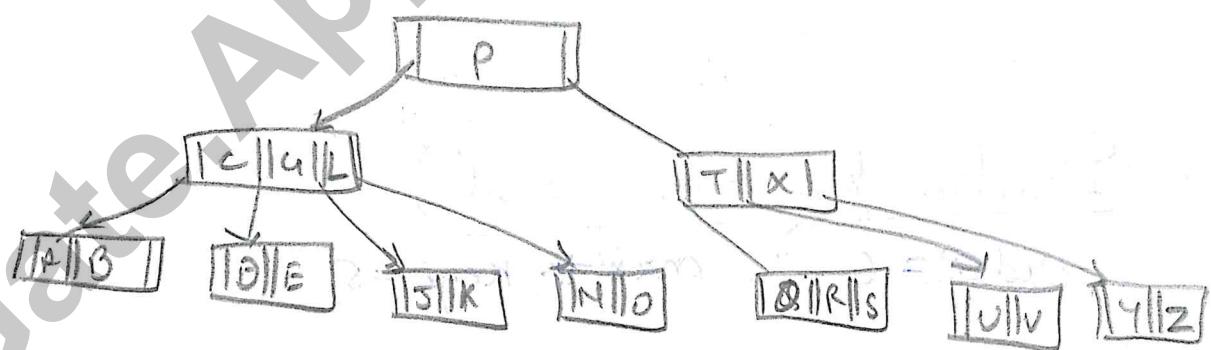
Phone: +91 844-844-0102

APPLIED ROOTS will simply remove ~~F~~ F from the leaf node. The resultant tree will be:



② Delete M :

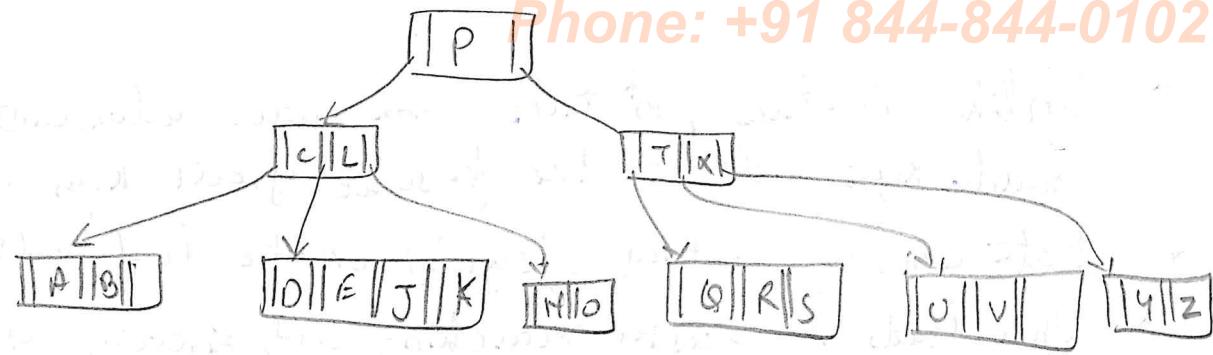
\because M is an internal node. The node before M has all values $< M$ and the node after M has all values $> M$. \therefore most logical value which should replace M is L.



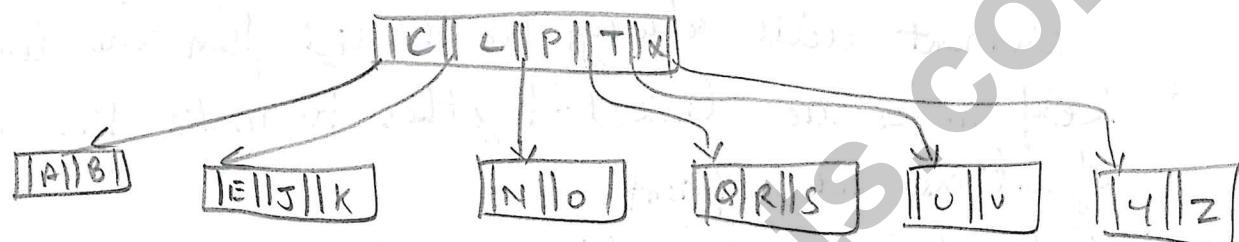
③ Delete G :

Since we can have max of keys = 5 \therefore either we can combine the node after deleting G or could do the same thing as we did in step ②)

Mail: gatecse@appliedroots.com



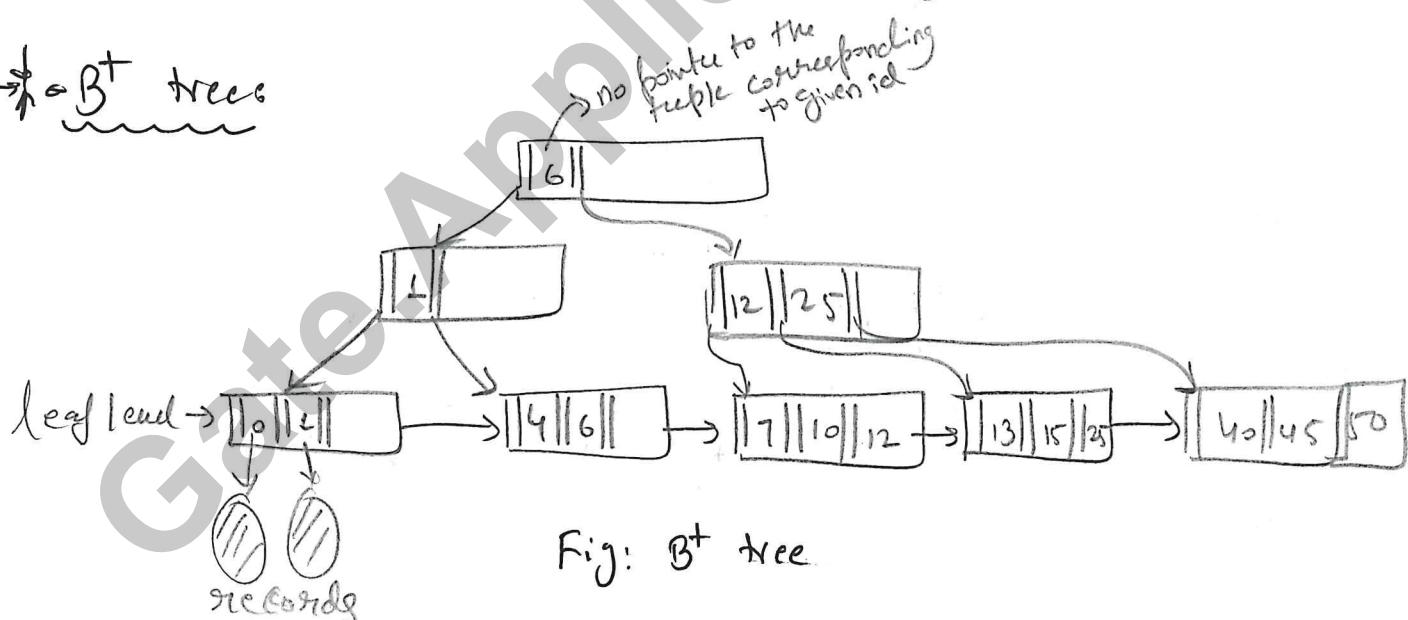
⑨ Delete D:



Note:-

While inserting or deleting ~~node~~ key in B-tree or B⁺ tree there are lots of cases and boundary cases to be handled.

~~B⁺~~ trees



- * Every node is a block
- * Root node has ≥ 2 pointers or children
- * Internal nodes have $\geq \lceil \frac{m}{2} \rceil$ pointers or children.
- * All leaves must be at some level

- * Unlike B-trees, B⁺ trees can have redundant search keys (due to the presence of all keys at leaf)
- * Data can only be stored (Records) on the leaf nodes.
- * This leads to faster searching comparatively as data can only be found on the leaf nodes.
- * Deletion will never be a complicated process since element will always be deleted from the leaf nodes.
- * Leaf nodes are linked together to make the search operations more efficient.
- * Order of leaf and non-leaf is different.

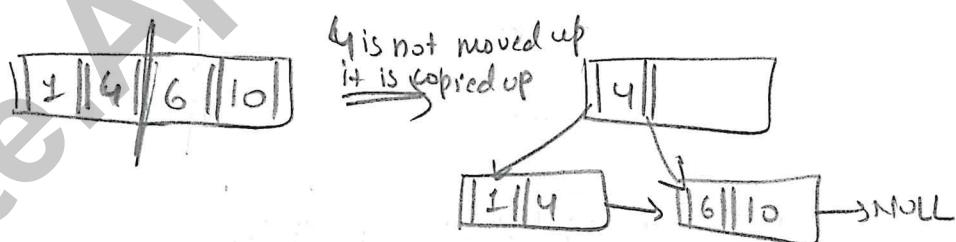
Insert into B⁺ trees:

$$\text{Order} = 4 \Rightarrow \# \text{Keys} = 3$$

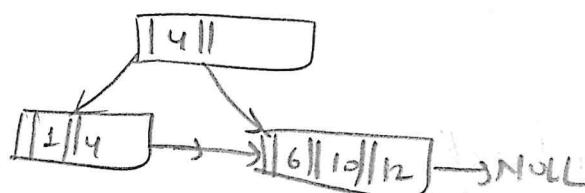
① Insert 1, 4, 6:



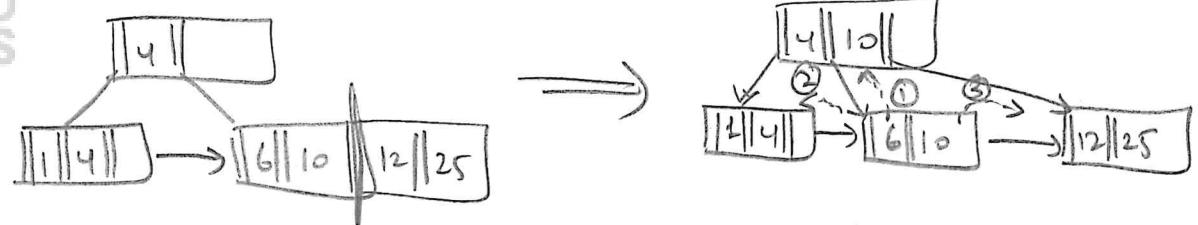
② Insert 10:



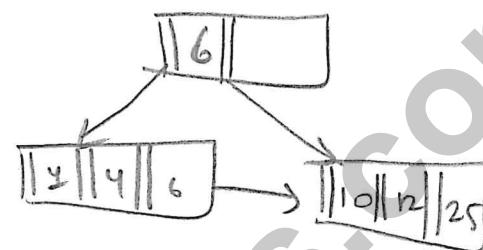
③ Insert 12:



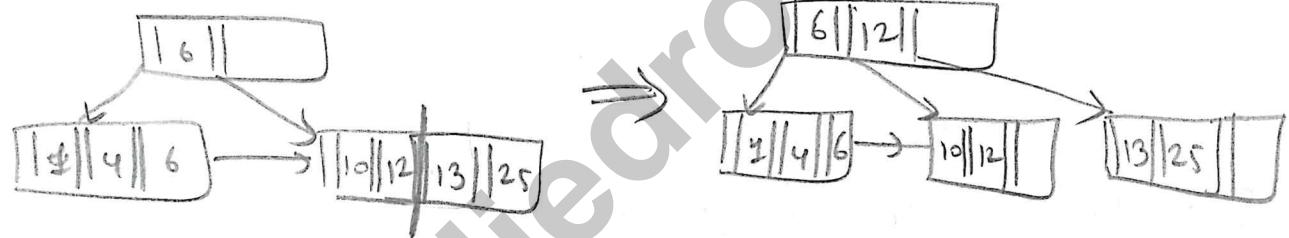
(4) Insert 25:



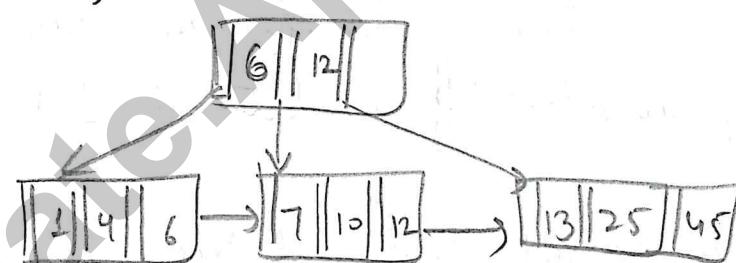
↓ Redistribution



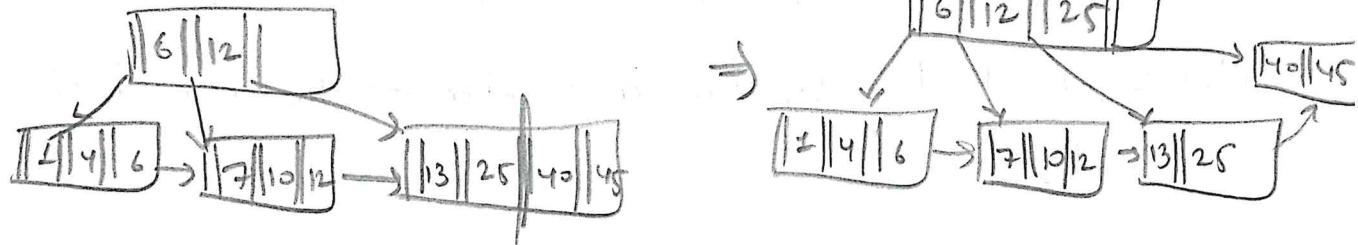
(5) Insert 13:



(6) Insert 7, 45:

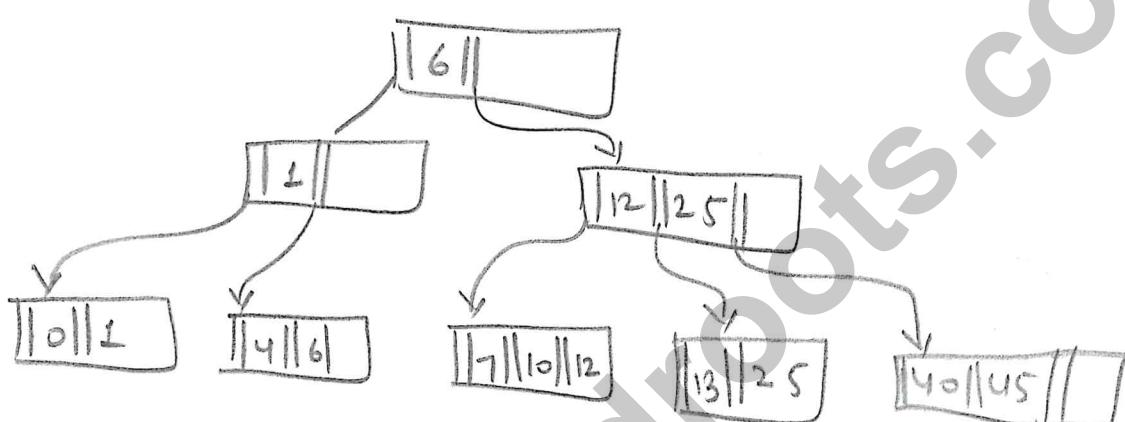
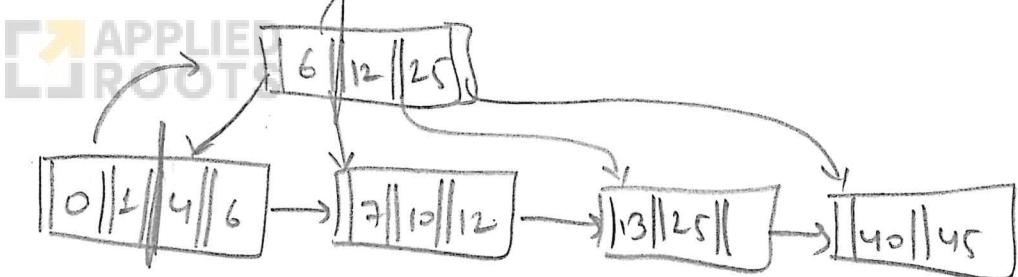


(7) Insert 40:



⑧ Insert 0 :

Phone: +91 844-844-0102



19.3 More solved problems:

(Q1) Which one of the following statement is NOT correct about the B⁺ tree data structure used for ~~for~~ creating an index of a relational database table?

- (A) B⁺ Tree is a height balanced tree
- (B) Non-leaf nodes have pointers to data records
- (C) Key values in each node are kept in sorted order
- (D) Each leaf node has a pointer to the next leaf nodes.

\Rightarrow (A) It is correct

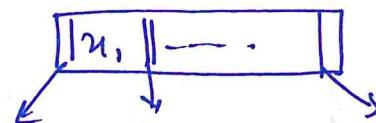
Phone: +91 844-844-0102

(B) In B^+ trees we point to leaf nodes. It is incorrect.

Hence, correct option is (B)

(Q2) In a B^+ tree, if the search key value is 8 bytes long, the block size is 512 bytes and the block pointer is 2 bytes, then the maximum order of the B^+ tree is _____

\Rightarrow



Let there be m keys then order will be $m+1$ pointers

$$BS = 512$$

$$P_b = 2$$

$$K = 8$$

$$\Rightarrow 8m + 2(m+1) \leq 512$$

$$\Rightarrow 10m + 2 \leq 512$$

$$\Rightarrow 10m \leq 510$$

$$\Rightarrow m = \lfloor \frac{510}{10} \rfloor = 51$$

(Q3) B^+ trees are considered BALANCED because

(A) The length of the paths from the root to all leaf nodes are all equal

(B) The length of the paths from the root to all leaf nodes differ from each other by at most 1.

- (C) The number of children of any two non-leaf sibling nodes differ by at most 1
- (D) The number of records in any two leaf nodes differ by at most 1

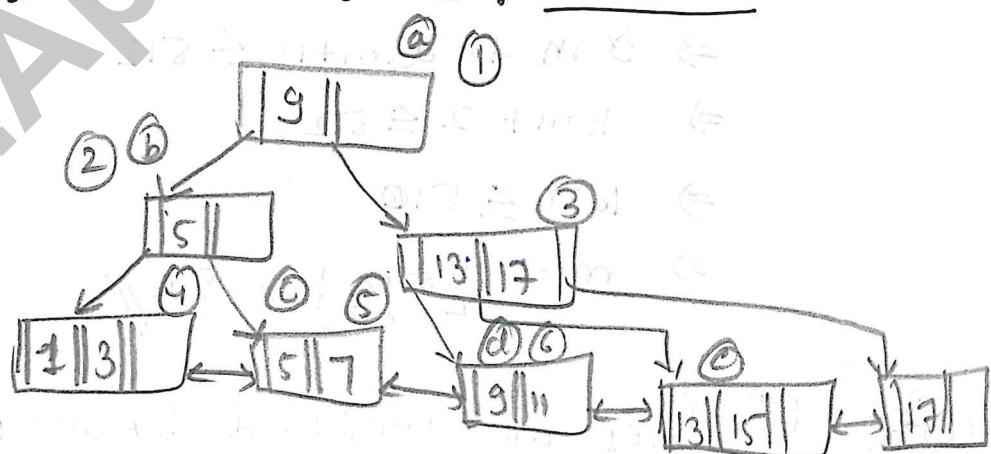
⇒ (A) is true
 (B) not true
 (C) not true
 (D) not true

∴ Correct option is (A)

- Q4) With reference to the B^+ tree index of order 3 shown below, the minimum no. of nodes (including the root node) that must be fetched in order to satisfy the following query:

"Get all records with a search key greater than or equal to 7 and less than 15" is:

- (A) 4
- (B) 5
- (C) 6
- (D) 7



⇒ We want minimum no. of nodes fetched while searching ≥ 7 and < 15 . On following search based approach will result in total 6 disk-read. (Here, ①, ②, ③, ④, ⑤, ⑥ are disk reads)