

2016

Wk-41/Day (278-088)

04  
Tuesday / 10

## Database Management System

Data - Raw and isolated facts about an entity (recorded) like text, audio, image, etc.

Information - Processed, meaningful and usable data

Database - Collection of similar / related data

DBMS - Software used to create, manipulate and delete database.

### Disadvantages of file system:

- i Data redundancy
- ii Data inconsistency
- iii Difficulty in accessing data
- iv Data isolation
- v Security problem
- vi Atomicity Problem
- vii Concurrent-access anomalies
- viii Integrity Problem

OCTOBER	
M	31
T	3 10 17 24
W	4 11 18 25
T	5 12 19 26
F	6 13 20 27
S	7 14 21 28
S	8 15 22 29
S	9 16 23 30

Phone/Email/Notes

Notes

P.T.O →

Comparison between two different types of databases :

OLAP (Online Analytical Processing)

OLTP (Online Transaction Processing)

- |                                  |  |
|----------------------------------|--|
| i) Historical data               | i) Current data                              |
| ii) Subject oriented             | ii) Application oriented                     |
| iii) Decision making             | iii) Day to day operations                   |
| iv) Larger size of data          | iv) Smaller size of data                     |
| v) Only read operations required | vi) Read and write, both operations required |

Entity Relationship Diagram : It was introduced by Dr. Peter Chen in 1976. It is a non-technical design method working on conceptual level based on the perception of real world.

Phone/Email/Notes

Notes

• Consists of collections basic objects called entities and relationships among these objects and attributes which defines their properties.

SEPTEMBER											
M	5	12	19	26							
T	6	13	20	27							
W	7	14	21	28							
T	1	8	15	22	29						
F	2	9	16	23	30						
S	3	10	17	24							
S	4	11	18	25							

- Free from ambiguities, and provides a standard and logical way of visualizing data.
- Basically, it is a diagrammatic representation and it is easy to understand even by non-technical user.

**Entity:** An entity is a thing or an object in the real world that is distinguishable from other objects/things based on the values of the attributes, it ~~possesses~~ possesses.

**Types of Entities :**

- ① **Tangible** - Entities which physically exist in real world. E.g., car, pen, etc.
- ② **Intangible** - Entities which exists logically, E.g., bank account, etc.

**Entity Set :** Collection of same type of entities i.e. those share the same properties or attributes

OCTOBER						
M	31	3	10	17	24	
T	1	6	11	18	25	
W	2	7	12	19	26	
TH	3	8	13	20	27	
F	4	9	14	21	28	
S	5	10	15	22	29	
SU	6	11	18	25		

Phone/Email/Notes

Notes

P.T.O →

Entity cannot be represented in an ER diagram as it is an instance / data. It can be represented in a relational model by row / tuple / record.

Entity set is represented by rectangle in ER diagram and is represented by table in a relational model.

Attributes: These are the units that describe the characteristics of entities.

- For each attribute, there is a set of permitted values called domains.
- In ER diagram, represented by ellipse or oval, while in relational model by a separate column.

Types:- ① Simple-Composite : ~~It~~, cannot be divided further & is simply represented by an oval.

② Single multivalued : It can have only one value at an instance of time.

Composite can be further divided in simple attribute where one oval can be connected with another.

	SEPTEMBER						
M	5	12	19	26	3	10	17
T	6	13	20	27	4	11	18
W	7	14	21	28			
T	1	8	15	22	29		
F	2	9	16	23	30		
S	3	10	17	24			
S	4	11	18	25			

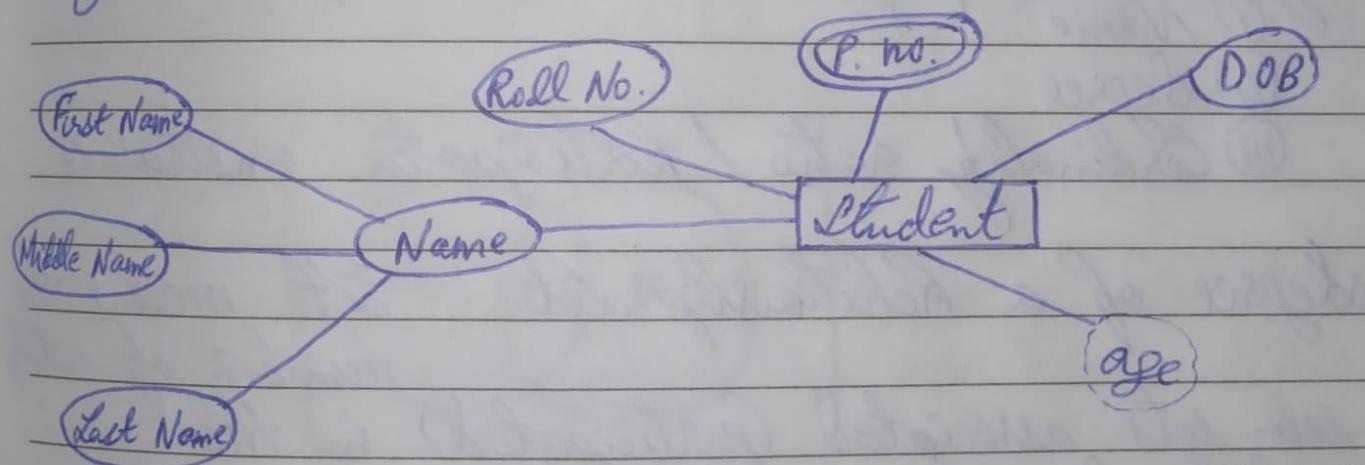
Single

~~Single~~ - Multivalued : Single can have only one value at an instance of time.

Multi Multivalued can have more than one value at an instance of time.

⑩ Store-Derived : stored is how value is stored in database.

Derived is how value can be computed in runtime using stored attribute



Sunday 09

Relationship : It is an association between two or more entities of same or different different entity set.

OCTOBER						
M	31	3	10	17	24	
T		4	11	18	25	
W		5	12	19	26	
T		6	13	20	27	
F		7				
S	1	8	14	21	28	
S	2	9	15	22	29	
		16	23	30		

It has no representation in ER diagram as it is an instance or a data.

In relational model, represented using

10/ Monday,

a row in a table.

Relationship type / set : A set of similar types of relationship.

Represented using a set diamond in ER diagram.  
In relational model, either by a separate table or by a separate column. (foreign key)

Every relationship type has three components -

- i) Name
- ii) Degree
- iii) Cardinality ratio / participation constraints

Degree of a Relationship set : It means the number of entity sets associated (participated) in the relationship set.

Most of the relationship sets in the E-R diagram are binary. Occasionally, however relationship sets involve more than two entity sets.

Phone/Email/Notes

Min. degree of relationship set = 1

Notes

	SEPTEMBER						
M	5	12	19	26			
T	6	13	20	27			
W	7	14	21	28			
T	1	8	15	22	29		
F	2	9	16	23	30		
S	3	10	17	24			
S	4	11	18	25			

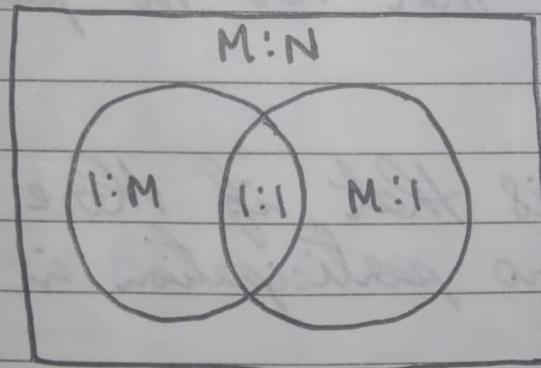
P.T.O →

# Mapping Cardinalities / Cardinality Ratio :

It expresses the number of entities to which other entity can be related via a relationship. It can be used in describing relationship set of any degree but is most useful in binary relationship.

1:1 (One to one)      M:1 (Many to one)

1:M (One to many)      M:N (Many to many)



Participation and Constraints : Specifies whether the instance of an entity depends on its relation to another entity via a relationship type.

These constraints describe the max and min no. of relationship instances that each entity can/must participate in.

OCTOBER	
M	31
T	3 10 17 24
W	4 11 18 25
T	5 12 19 26
F	8 13 20 27
S	7 14 21 28
S	1 8 15 22 29
S	2 9 16 23 30

Phone/Email/Notes

Notes

**Max Cardinality** : It defines the max number of times, an entity occurrence participates in a relationship.

**Min Cardinality** : It defines the min no. of times, an entity occurrence participates in a relationship.

**Notation** : (min c, max c)

Partial participation is that if there exists at least one entity that has no participation in the relationship.

Total participation is that if there exists no entity which has no participation in the relationship.

**Functional Dependencies** :

$\beta: \alpha \rightarrow \beta$  Here, using the value of  $\alpha$ , we can map it in the table and search for the value of  $\beta$ .

	$\alpha$	$\beta$
$t_1 \rightarrow$	a	1
$t_2 \rightarrow$	a	1
	c	2
	d	3

Phone/Email/Notes

Notes

SEPTEMBER						
M	5	12	19	26		
T	6	13	20	27		
W	7	14	21	28		
T	1	8	15	22	29	
F	2	9	16	23	30	
S	3	10	17	24		
S	4	11	18	25		

If  $t_1[\alpha] = t_2[\alpha] \Rightarrow t_1[\beta] = t_2[\beta]$

$\alpha$  is called determinant &  $\beta$  is called dependent.

$$f: \alpha \rightarrow \beta$$

Trivial  
 $\beta \subseteq \alpha$

Non-trivial  
 $\beta \not\subseteq \alpha$

Q:

R

A	B	C	D	E
a	2	3	4	5
b	a	3	4	5
a	2	3	6	5
a	2	3	6	6

- @ A  $\rightarrow$  BC
- @ DE  $\rightarrow$  C
- @ C  $\rightarrow$  DE
- @ BC  $\rightarrow$  A

- Sol:
- @ A  $\rightarrow$  BC is correct
  - @ DE  $\rightarrow$  C is correct
  - @ BC  $\rightarrow$  A is correct

We check if on same value of  $\alpha$ , we get two different values of  $\beta$ . If so, then the functional dependency is not valid. Otherwise, the functional dependency is valid.

OCTOBER				
M	31	10	17	24
T	4	11	18	25
W	5	12	19	26
T	6	13	20	27
F	7	14	21	28
S	8	15	22	29
S	9	16	23	30

Phone/Email/Notes

- \* If all the values of  $\alpha$  are unique f.d. is always valid.
- \* If all the values of  $\beta$  are same, f.d. is always valid.

Attribute Closure / Closure on Attribute set / Closure of Attribute set :

Attribute Closure of an attribute set 'A' can be defined as a set of attributes which can be functionally determined from it. Denoted by  $F^+$ .

$R(ABCDEF) : A \rightarrow B, C \rightarrow DE, AC \rightarrow F,$   
 $D \rightarrow AF, E \rightarrow CF$

From above,  $(D)^+ = ABDF$ ,  $(DE)^+ = ABCDEF$

Armstrong Axioms/Rules : These hold true for relational databases & can be used to generate closure set.

Primary Rules - Reflexivity : If  $Y \subseteq X$ ,  
 Then  $X \rightarrow Y$

Augmentation : If  $X \rightarrow Y$ , then  $XZ \rightarrow YZ$

Transitivity : If  $X \rightarrow Y$  &  $Y \rightarrow Z$ , then  $X \rightarrow Z$

Phone/Email/Notes

Notes

Secondary Rules - Union : If  $X \rightarrow Y$   
 &  $Y \rightarrow Z$ ,  
 Then  $X \rightarrow YZ$

SEPTEMBER											
M							5	12	19	26	3
T							6	13	20	27	4
W							7	14	21	28	5
T							1	8	15	22	6
F							2	9	16	23	7
S							3	10	17	24	8
S							4	11	18	25	9

Decomposition: If  $X \rightarrow YZ$ , then  $X \rightarrow Y$  &  $X \rightarrow Z$

Recom. Transitivity: If  $X \rightarrow Y$  &  $Y \rightarrow Z$ ,  
then  $X \rightarrow Z$

Composition: If  $X \rightarrow Y$  &  $Z \rightarrow W$ , then  
 $XZ \rightarrow YW$

Closure Set of Attributes:

~~R(ABC)  $A \rightarrow B$   $B \rightarrow C$~~

~~Let  $F = F_1 + F_2$   
 $A \rightarrow B$      $A \rightarrow C$~~      $\therefore A \nrightarrow BC$

Equivalence of Functional Dependency:

Sunday 16

~~R(ACDEH)~~     $F: A \rightarrow C$      $g: A \rightarrow CD$   
 $Ac \rightarrow D$                    $E \rightarrow AH$   
 $E \rightarrow AD$   
 $B \rightarrow H$

OCTOBER						
M	31	3	10	17	24	
T	4	11	18	25		
W	5	12	19	26		
F	6	13	20	27		
S	7	14	21	28		
S	8	15	22	29		
S	9	16	23	30		

$\textcircled{Q} F \subseteq g$      $\textcircled{Q} F \supseteq g$   
 $\textcircled{Q} F = g$      $\textcircled{Q} F \neq g$

Phone/Email/Notes

Notes

sol:

P.T.O →

Finding closure of LHS of  $F$  using a

$$(A)^+ = ACD$$

$$(AC)^+ = ACD$$

$$(E)^+ = EAHCD$$

Comparing this with the functional dependencies of  $F$

$\therefore$  All the dependencies of  $a$  can be done by  $F$

$$\therefore F \subseteq a$$

Finding closure of L.H.S of  $a$  using  $F$

$$(A)^+ = ACD$$

$$(E)^+ = EADHC$$

Comparing this with the f.d. of  $a$

$\therefore$  All the dependencies of  $F$  can be done by  $a$

$$\therefore a \subseteq F$$

$$\therefore F = a$$

Irreducible set of functional dependency (Canonical form):

Phone/Email/Notes

$R(WXYZ)$   $X \rightarrow W$

Notes

$WZ \rightarrow XY$

$Y \rightarrow WXZ$

SEPTEMBER						
M	5	12	19	26		
T	6	13	20	27		
W	7	14	21	28		
T	1	8	15	22		
F	2	9	16	23		
S	3	10	17	24		
S	4	11	18	25		

Decomposing these f.d., we have

$\checkmark X \rightarrow W$	$(X)^+ = XW$	$(X)^+ = X$
$XWZ \rightarrow X$	$(WZ)^+ = WZXZ$	$(WZ)^+ = WZYX$
$\checkmark WZ \rightarrow Y$		$(WZ)^+ = WZ$
$XY \rightarrow W$	$(Y)^+ = YWXZ$	$(Y)^+ = YXZW$
$\checkmark Y \rightarrow X$		$(Y)^+ = YZ$
$\checkmark Y \rightarrow Z$		$(Y)^+ = YXW$

After filtering, we have

$$X \rightarrow W, WZ \rightarrow Y, Y \rightarrow X, Y \rightarrow Z$$

Checking the 'x' values of f.d.,

$$(WZ)^+ = WZYX$$

$$(W)^+ = W \quad (Z)^+ = Z$$

Canonical Canonical form :  $X \rightarrow W$   
 $WZ \rightarrow Y$   
 $Y \rightarrow XZ$

\* There can be different minimal forms for a same set of f.d.

OCTOBER	
M	31
3	10
17	24
W	4
11	18
25	
T	5
12	19
26	
F	6
13	20
27	
S	7
14	21
28	
O	8
15	22
29	
Z	2
9	16
23	30

Keys : When we find the closure set of key, we get the complete relation

$$\text{Key}^+ = R$$

It is a set of attributes which can uniquely identify a row or a tuple in a relational table.

i)  $R(ABCD)$ ;  $A \rightarrow BC$

Here, A is not key as  $(A)^+ = ABC \neq R$

ii)  $R(ABCD)$ ;  $ABC \rightarrow D$ ,  $AB \rightarrow CD$ ,  $A \rightarrow BCD$

Here,  $(ABC)^+ = (AB)^+ = (A)^+ = R$ . Thus, all are keys

Super Key: It is a set of attributes with the help of which you can uniquely identify ~~a~~ all the remaining attributes. This is the most basic key.

Candidate Key: A super key is candidate key whose no proper subset is a super key again.

iii)  $R(ABCD)$ ;  $B \rightarrow ACD$ ,  $ACD \rightarrow B$

Phone/Email/Notes

Notes Here,  $(B)^+ = (ACD)^+ = R$

Both B and ACD are candidate keys.

SEPTEMBER	
M	5
T	6
W	7
T	1
F	2
S	3
S	4
	12
	13
	14
	15
	16
	17
	18
	19
	20
	21
	22
	23
	24
	25
	26
	27
	28
	29
	30
	31

Primary Key: It is that candidate key which is selected by a database administrator as a primary key way to identify tuples.

⑩  $R(ABCD)$ ;  $AB \rightarrow C$ ,  $C \rightarrow BD$ ,  $D \rightarrow A$

$(AB)^T = ABCD = R$ ,  $(C)^T = CBDA = R$ ,  $(D)^T = DA$

Here, AB, C are <sup>Super</sup>candidate keys and candidate keys.

Normalization:

① First Normal Form - Every cell should have a single value.  
(1NF)

If there is a multivalued attribute present, we ~~dep~~ duplicate all the other attributes and convert it into 1NF.

In general, for every multivalued attribute, we make a separate table.

OCTOBER	
M	31
T	3 10 17 24
W	4 11 18 25
T	5 12 19 26
F	6 13 20 27
S	7 14 21 28
S	1 8 15 22 29
S	2 9 16 23 30

Here, the order of rows and the order of columns ~~columns~~ columns Notes  
Notes  
This is irrelevant. Also, every column should have a unique name.

(ii) Second Normal Form -

$R(A \downarrow B \downarrow C \downarrow D)$

$$(AB)^+ = ABCD \\ = R$$

Here, A, B are prime attribute and others are non prime attribute.

Here, C is only dep. When a non-prime attribute instead of depending on the entire candidate key, is dependent on a part of it is called partial dependent dependency.

Thus, for a table to be in 2NF, it should be in 1NF and it shouldn't have any partial dep. dependency.

We decompose  $R(ABCD)$  into two tables

$R_1(A \downarrow B \downarrow D)$ ,  $R_2(B \downarrow C)$

One table is for the complete candidate key and all the complete attributes having no partial dependencies.

Notes      Another table is for the partial dependency.

This, we obtain 2NF

SEPTEMBER						
M	5	12	19	26		
T	6	13	20	27		
W	7	14	21	28		
T	1	8	15	22	29	
F	2	9	16	23	30	
S	3	10	17	24		
S	4	11	18	25		

Third Normal Form - A functional dependency from  $\alpha \rightarrow \beta$  is called transitive if  $\alpha, \beta$  are non prime.

$R(ABC)$ ;  $A \rightarrow B, B \rightarrow C$

A relation R is in the 3NF if it is in 2NF and there should not be transitive dependency.

Also if every dependency from  $\alpha \rightarrow \beta$  where  $\alpha$  is a super key or  $\beta$  is a prime attribute, then it is in 3NF.

Boyce-Codd Normal Form - It deals with  $\alpha \rightarrow \beta$  p/n.p  
(BCNF)

$R(ABC)$ ;  $(AB)^+ = ABC = R$        $AB \rightarrow C$       Sunday 23  
                 $(AC)^+ = ABC = R$        $C \rightarrow B$

Both the f.d. are neither partial nor transitive. Thus, the table is in 2NF and 3NF.

for  $\alpha$  for the table to be  
in BCNF,  $\alpha$  should be super key.

OCTOBER						
M	31	3	10	17	24	
T	4	11	18	25		
W	5	12	19	26		
T	6	13	20	27		
F	7	14	21	28		
S	8	15	22	29		
S	9	16	23	30		

Phone/Email/Notes

Notes



Normalising the table,

$R(ABC)$

$R_1(CB)$

$R_2(AC)$

Ways to check normal form :

- i) For every dependency from  $\alpha \rightarrow \beta$ , if  $\alpha$  is super key, then it is BCNF.
- ii) For every dependency from  $\alpha \rightarrow \beta$ , if  $\alpha$  is candidate key,  $\beta$  is non prime, then it is 3NF.
- iii) For every dependency from  $\alpha \rightarrow \beta$ , if  $\beta$  is a non prime, then it is 2NF.
- iv) For every dependency from  $\alpha \rightarrow \beta$ , if  $\alpha$  is C.K,  $\beta$  is prime, then it is 1NF.

Lossless join decomposition / non-additive :

This property guarantees that any extra or less row/tuple generation problem does not occur after decomposition.

- It is a mandatory property & must always hold good.

SEPTEMBER						
M	5	12	19	26	3	10
T	6	13	20	27	4	11
W	7	14	21	28	5	12
T	8	15	22	29	6	13
F	9	16	23	30	7	14
S	10	17	24	31	8	15
S	11	18	25	1	9	16

- If a relation  $R$  is decomposed into two relations  $R_1$  &  $R_2$ , then it will be lossless if and only if
  - $\text{attr}(R_1) \cup \text{attr}(R_2) = \text{attr}(R)$
  - $\text{attr}(R_1) \cap \text{attr}(R_2) \neq \emptyset$
  - $\text{attr}(R_1) \cap \text{attr}(R_2) \rightarrow \text{attr attr}(R_1)$   
or  
 $\text{attr}(R_1) \cap \text{attr}(R_2) \rightarrow \text{attr}(R_2)$

Dependency Preservation: If a table  $R$  having f.d. set  $F$  is decomposed into two tables  $R_1$  and  $R_2$  having f.d. set  $F_1$  &  $F_2$  then  $F_1 \subseteq F^+$ ,  $F_2 \subseteq F^+$  i.e.  $(F_1 \cup F_2) = F^+$ . This means the decomposition is dependency preserving.

$R(ABCD)$ ;  $AB \rightarrow CD$ ,  $D \rightarrow A$

$R_1(AD)$

$F_1: D \rightarrow A \checkmark$   
 $A \rightarrow D \times$

$R_2(BCD)$

$F_2: D \rightarrow A \times$   
 $CD \rightarrow CDA \times$   
 $BD \rightarrow BDAC \checkmark$

OCTOBER	
M	31
T	8 10 17 24
W	4 11 18 25
T	5 12 19 26
F	6 13 20 27
S	7 14 21 28
S	8 15 22 29
	9 16 23 30

Phone/Email/Notes

Notes

The dependency  $AB \rightarrow CD$  is lost.  
Dependency preservation is not compulsory while table decomposition.

\* Those attributes which don't have any incoming edge are called essential attributes.

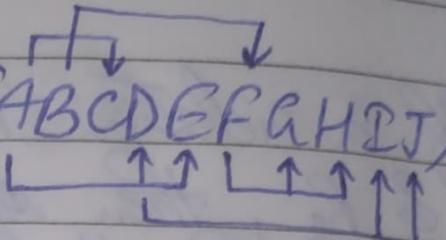
28

10

Friday

Wk-44/Day (302-064)

Normalisation of a Table:  $R(ABCDEF\overset{\uparrow}{G}\overset{\uparrow}{H}\overset{\uparrow}{I}\overset{\uparrow}{J})$



$$(AB)^+ = ABDE \overset{\uparrow}{C} \overset{\uparrow}{F} \overset{\uparrow}{G} \overset{\uparrow}{H} \overset{\uparrow}{I} \overset{\uparrow}{J} \rightarrow \text{c.k.}$$

The current  $R$  is in 1NF. We normalise it further.

$$R_1(\overset{\uparrow}{ABC}), R_2(\overset{\uparrow}{ADE}\overset{\uparrow}{IJ}), R_3(\overset{\uparrow}{BF}\overset{\uparrow}{GH})$$

Current schema is in 2NF. We normalise it further.

$$R_{11}(\overset{\uparrow}{ABC}), R_{21}(\overset{\uparrow}{ADE}), R_{22}(\overset{\uparrow}{D}\overset{\uparrow}{IJ}), R_{31}(\overset{\uparrow}{BF}), R_{32}(\overset{\uparrow}{FG})$$

Current schema is in 3NF & BCNF as well.

# Database Management System

## File Structures :

### Sorted file

- 1) Can be sorted only according to one alternate attribute (Search key).
- 2) Searching is fast (Binary search).
- 3) Insertion and deletion can be difficult.

### Unsorted file

- 1) Random order, any record can be placed anywhere.
- 2) Search is slow (in linear search).
- 3) Insertion and deletion can be easy.

## Spanned / Unspanned Mapping :

Spanned is when a table is broken into parts while storing if the block cannot store a complete table. costly in time.

Unspanned is when a table is stored in a block if and only if the block can store a complete table. costly in space.

Indexing : We make a file called as the index file. We store two columns, Search Attribute and Block pointer.

\* The block pointer contains the base address of the whole of that block which contains the ~~the~~ search attribute.

If selected records get stored in the index file, then it is known as sparse indexing.

If all values records get stored in the index file, then it is

known as dense indexing.

## Types of Indexing:

We can make multiple index files depending upon the size of the index file. It is known as multi-level indexing.

- ① Primary Indexing - Here, the search attribute is the primary key of the table. It is done when the main file is sorted.
- ② Clustered Indexing - Here, the search attribute is not a key of the table.
- ③ Secondary Indexing - Here, the search attribute is either key or non-key and the main file is not sorted.

\* Index file will always be sorted.

Primary Indexing: ① The main file is sorted.

② Primary key is used as search attribute.

③ It is an example of sparse indexing.

④ No. of entries in the index file = No. of blocks acquired by the main file.

⑤ No. of access required =  $\log_2 n + 1$

Let us assume a block size of 1024 MB = 1024B

$$\text{Blocking factor} = \left\lceil \frac{\text{Block size}}{\text{Record size}} \right\rceil = \left\lceil \frac{1024B}{15B} \right\rceil = 68.26$$

$$= [68.26] = 68$$

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

$$\therefore \text{No. of block access required} = \lceil \log_2 15 \rceil + 1 = 6 + 1 = 7$$

Clustered Indexing: (i) The main file is sorted (on a non-key attribute).

(ii) There will be only one ~~entry~~ entry for each unique value of the non-key attribute.

(iii) If the number of blocks acquired by index file is  $n$ , then block access required will be  $\geq \log_2 n + 1$ .

Secondary Indexing: (i) The main file is unsorted.

(ii) It can be done on key as well as non-key attributes.

(iii) It is called secondary because normally one indexing is already done.

(iv) It is an example of dense indexing, no. of entries in index file = no. of entries in main file.

$$\text{(v) No. of block access} = \lceil \log_2 n \rceil + 1$$

Let us consider a block size of 1024B. The record size in main file is 100B and in index file is 15B. The no. of records in main file are 30000.

$$\text{Blocking factor} = \left\lceil \frac{\text{Block size}}{\text{Record size}} \right\rceil = \left\lceil \frac{1024B}{15B} \right\rceil = 68$$

$$\text{No. of blocks for index file} = \left\lceil \frac{\text{No. of records for index file}}{\text{Blocking factor}} \right\rceil = \left\lceil \frac{30000}{68} \right\rceil = 442$$

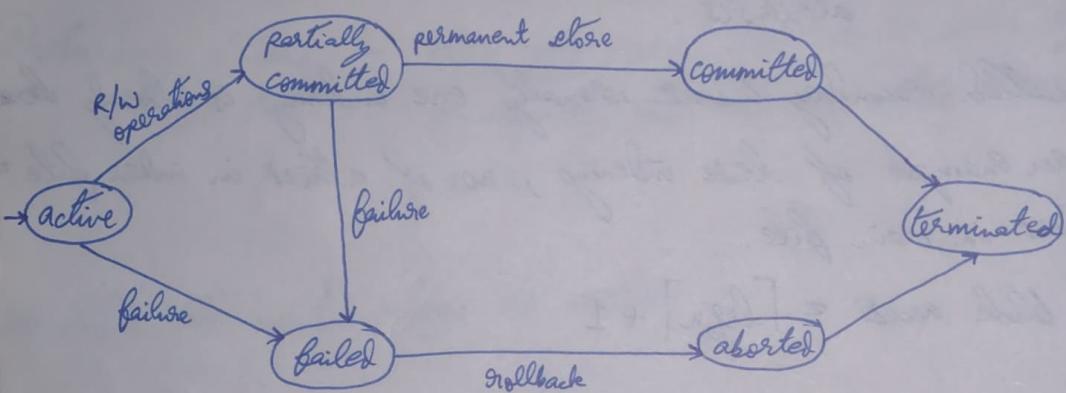
$$\therefore \text{No. of block access required} = \lceil \log_2 n \rceil + 1 = 9 + 1 = 10$$

Transactions: It is a set of SQL instructions which perform a logical unit of work. They are atomic in nature.

ACID Properties:

- i) **Atomicity** :- It means the transactions are atomic in nature.
- ii) **Consistency**: After the transaction is completed, the DB should be consistent.
- iii) **Isolation**: Each transaction should be logically isolated of the other.
- iv) **Durability**: The transaction carried out should be durable in nature.

Transaction States:



- \* A committed transaction cannot be rolled back.
- \* Concurrency control component takes care of isolation.
- \* Recovery management component takes care of durability.

Advantages of Concurrency:

- i) **Waiting time** ↓ - After a process is ready, the time it waits until execution starts is known as waiting time.
- ii) **Response time** ↓ - The time taken by the first response of the CPU after the process is ready is known as response response time.
- iii) **Resource Utilization** ↑ - It is the utilization of the resources.

for carrying out the concurrent tasks.

(iv) Efficiency ↑ - It determines how much resources ~~yes~~ have been used for carrying out the concurrent tasks.

$$\text{Efficiency} = \frac{\text{No. of tasks}}{\text{No. of resources}} \times 100\%$$

t. Dirty Read problem : In case of concurrent transactions, when one transaction reads data from the local buffer of another transaction, and commits it, then it is known as dirty read. The problem occurs when there is a failure in the other transaction. This leads to rollback of the transaction and the value data read and committed by the first transaction doesn't exist in the database anymore. This leads to inconsistency.

\* To solve dirty read problem, the first transaction should only commit after the second transaction has committed. This allows the first transaction to rollback if necessary.

Unrepeatable Read problem : In case of concurrent transactions, when there are multiple reads of data by some transaction and it also gives different output each time, then it is known as unrepeatable read problem.

Phantom Read problem : In case of concurrent transactions, when there are multiple reads of data by some transaction and it is unable to read data at some point, then it is known as phantom read problem.

Lost Update problem (Ulrite-Write conflict) : In case of concurrent transactions, when there

are multiple transactions writing to the same database, the data updated by the previous transactions get lost as the data gets overwritten by other transactions and finally committed. This is known as lost update problem.

Schedule : It is used to schedule concurrent transactions.

Serial Schedule - Transactions are executed one after the other.

Non serial Schedule - Transactions are switched among one another and are executed.

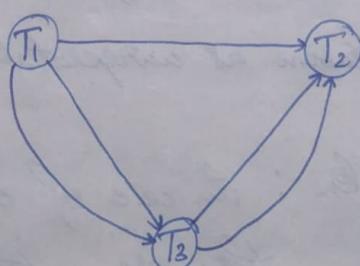
\* For  $n$  transactions,  $n!$  serial schedules are possible.

\*  $n = n_1 + n_2 + \dots + n_n$ ,  $\frac{(n_1 + n_2 + \dots + n_n)!}{n_1! \cdot n_2! \cdot \dots \cdot n_n!} = n!$  non serial schedules are possible.

Conflict Serializability : We use this to convert a non serial schedule to a serial schedule.

\* Conflict occurs when there are ~~RFW~~<sup>write instructions</sup> operations on the same data of the transactions concurrent transactions.

	T <sub>1</sub>	T <sub>2</sub>	T <sub>3</sub>	do:
$\times R(X)$				
$\times R(Y)$			$\times R(Y)$	
$\times W(Y)$				
$\times W(X)$			$\times W(X)$	
$\times R(X)$				
$\times W(X)$				



There are no cycles in this graph. Thus, the given schedule is conflict serializable.

There are no incoming edges to T<sub>1</sub>. Thus, T<sub>1</sub> will run first.

There are no outgoing edges from  $T_2$ . Thus,  $T_2$  will run last.  
∴ Order of serializability :  $T_1 \rightarrow T_3 \rightarrow T_2$

	$T_1$	$T_2$	$T_3$	Ans:	$T_1$	$T_2$
	$\times R(a)$					
		$\times R(b)$				
			$\times R(c)$			
			$\times W(c)$			
		$\times W(b)$				
			$\times W(a)$			

There are no cycles present in the graph. Thus, the schedule is conflict serializable.

∴ No. of serial schedules possible =  $3! = 6$

\*  $O(n^2)$  is the time complexity for finding out whether a schedule is conflict serializable or not.

View serializability: It is a superset of conflict serializability.  
A schedule is not conflict serializable but can be view serializable when there is atleast one blind write present in the schedule.

When the given non-serial schedule is view equivalent to any one of the  $n!$  serial schedules, then the given schedule is view serializable.

The conditions for View Equivalence :

- (i) The transaction which does the initial read on a data in first schedule, ~~should do same~~ that transaction only should do the initial read on that data in second schedule.
- (ii) The transaction which does the final write on a data in first

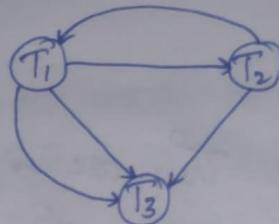
schedule, that transaction only should do the final write on that data in second schedule.

(iii) In first schedule, the intermediate reads done by a transaction should be done by the same transaction in the second schedule.

Q: )

$T_1$	$T_2$	$T_3$
$\times R(a)$		
	$\times W(a)$	
$\times W(a)$		
		$\times W(a)$

Ans:



The ~~go~~ There is a cycle present. So, the schedule is conflict serializable.

Now, there is blind write present in  $T_2$  and  $T_3$ . Thus, the schedule is not view serializable.

Taking a serial schedule,

$T_1$	$T_2$	$T_3$
$R(a)$		
	$W(a)$	
		$W(a)$
		$W(a)$

Comparing this with the given non serial schedule, we can observe that all the three conditions of  $\Rightarrow$  view equivalent is satisfied.

Thus, the given non-serial schedule is view-serializable.

Recoverable Schedule: These schedules can handle failures in transactions and ensure that the consistency is maintained.

\* If there no dirty read present in the schedule, then it is a recoverable schedule.

\* If there are dirty reads present in a schedule, then the transaction doing the dirty read should commit later to make it a recoverable schedule.

\* A schedule should always be recoverable.

### Cascadless Schedule

When multiple transactions rollback one after the other due to a failure, then it is known as cascading rollbacks.

When there are no cascading rollbacks occurring in a schedule, then it is known as cascadless schedule. To make a cascadless schedule, we remove dirty reads from the schedule.

Strict Schedule : In a strict schedule, there are no dirty reads by any other transaction as well as no write operation can be performed until the <sup>first</sup> original transaction commits the data.

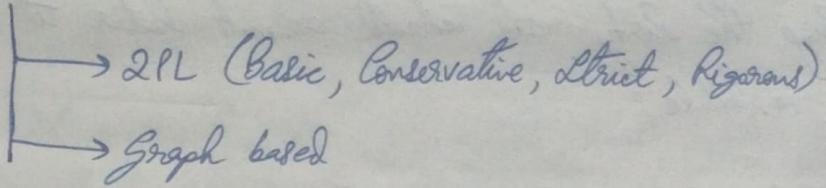
T <sub>i</sub>	S <sub>1</sub>		S <sub>2</sub>	
	T <sub>1</sub>	T <sub>2</sub>	T <sub>1</sub>	T <sub>2</sub>
	R(a)		R(a)	
	W(a)			R(b)
C ✓				
	Rf W(a)		W(a)	
	R(a)			W(b)
	C		C ✓	
				R(a)
				C

### Concurrency Control Technique

The main problem is different transactions trying to access data at the same time.

For that we use, ① Time stamping Protocol

## ii) Lock based protocol



## iii) Validation protocol

Time stamping protocol : The basic idea is to decide the order between the transactions before that enters into the system so that in case of conflict during execution, we can resolve the conflict using ordering.

The reason we call it time stamp and not just stamp because for stamping we use the value of the clock as it is always unique and never repeat itself.

Two ways of time stamping -

i) Time stamp with transaction : With each transaction  $t_i$ , we associate a time stamp denoted by  $T.S.(T_i)$ , it is the value of the system clock when a transaction enters into the system so if a new transaction  $T_j$ , enters after  $T_i$  then,  $T.S.(T_i) < T.S.(T_j)$  which will always be unique and will remain fixed throughout the execution. It also determines the serializability order if  $T.S.(T_i) < T.S.(T_j)$  then system ensures that the resultant conflict serializable schedule  $T_i$  will execute first, before  $T_j$  executes.

ii) Time stamp of each data item : For each data item  $Q$ , the protocol maintains two time stamps.

W-Timestamp ( $\theta$ ) is the largest time stamp of any transaction that enacted write ( $\theta$ ) successfully.

R-Timestamp ( $\theta$ ) is the largest time stamp of any transaction that enacted read ( $\theta$ ) successfully.

T<sub>i</sub> request for Read ( $\theta$ ) :

If  $T.S.(T_i) < W.T.S.(\theta)$ , means  $T_i$  needs to read a value of  $\theta$  that was already overwritten. Hence, request is rejected and  $T_i$  must rollback.

If  $T.S.(T_i) \geq W.T.S.(\theta)$ , operation can be allowed and R.T.S. ( $\theta$ ) will be  $\max(R.T.S.(\theta), T.S.(T_i))$ .

T<sub>i</sub> request for Write ( $\theta$ ) :

If  $T.S.(T_i) < R.T.S.(\theta)$ , means value of  $\theta$  that  $T_i$  is producing was not needed previously and the system assumed that the value would never be produced and hence rejected and rolled back.

If  $T.S.(T_i) < W.T.S.(\theta)$ ,  $T_i$  is attempting to write an obsolete value of  $\theta$ , thus  $T_i$  is rejected and rolled back otherwise it is allowed.  $W.T.S.(\theta) = \max(W.T.S.(\theta), T.S.(T_i))$ .

Properties of Time Stamping Protocol -

- i) It ensures conflict serializability.
- ii) It ensures view serializability.
- iii) Possibility of dirty reads and no restriction on ~~concurrent~~ commit, so ~~more~~ irrecoverable schedules and cascading rollbacks are possible.
- iv) No deadlock as we either allow or reject the transaction.
- v) Time-stamped schedules made by time stamping are a subset of conflict serializable schedules.

(iv) It may suffer from starvation and is relatively slow.

Thomas Write Rule : Modifies time stamping protocol to make some improvements which may generate those schedules that are V.S. but not C.S. and provides better concurrency.

It modifies time stamping protocol in ~~obsolete~~ a blind write case when  $T_i$  request  $w(Q)$  if  $T_i < T.S.(T_i) < W.T.S.(Q)$ .

Here,  $T_i$  attempts to write obsolete value of  $Q$  so rather than rolling back  $T_i$ , write operation operation is ignored.

Lock Based Protocol : To achieve consistency, isolation is the most important idea and locking is used to achieve isolation. First obtain a lock on a data item, then perform the desired desired operation and then unlock it.

Different modes of locks -

i) Shared mode : It is denoted by lock- $s(Q)$ . Transaction that can perform ~~as~~ read operation and any other transaction can also obtain the same lock on the same data item at the same time.

ii) Exclusive mode : It is denoted by lock- $x(Q)$ . Transaction can perform both R&W operation and any other transaction cannot be obtained either with shared or exclusive lock.

proper properties of lock based protocol -

i) If we do unlocking, inconsistency will arise and if we don't unlock, then concurrency will be poor.

ii) The transaction should follow a set of rules for locking and

unlocking of late item like 2-phase locking or graph based.

(iii) A schedule  $\sigma$  is legal under a protocol if it can be generated using the rules of the protocol.

Two-Phase Locking (2PL): (i) This protocol requires that each transaction in a schedulable schedule will be too two phased, Growing phase and Shrinking phase.

(ii) \* In growing growing phase, transaction can only obtain locks but cannot release any lock.

\* In shrinking phase, transaction can only release locks but cannot obtain any lock.

Transaction can perform both R&W operations, both in growing and shrinking ph phase. It ensures C.S. and V.S. and the order of serializability is the order in which transaction reaches the lock point. It may generate ~~in areas~~ irrecoverable schedules and cascading rollbacks and does not ensure freedom from deadlock.

Conservative/Static 2PL : There is no growing phase, transaction, first we will acquire all the locks ~~required~~ required and then directly will start from the lock point.

If all the locks are not available, then the transaction must release the lock acquired so far and wait. Shrinking phase works as usual and the transaction can unlock any data item.

\* C.S., V.S. and independent of deadlock but possibility of ~~recoverable~~ irrecoverable schedules and cascading rollbacks.

Rigorous 2PL : It is an improvement over 2PL protocol where we try

to ensure recoverability and cascadelessness.

- i) It requires that all the locks must be held until the transaction commits i.e. there is no shrinking phase in the life of a transaction.
- ii) Will ensure C.S., V.S., recoverability and cascadelessness.
- iii) Differs from deadlock and inefficiency.

Strict 2PL : It is an improvement over rigorous 2PL.

- i) In the shrinking phase, unlockings of one exclusive task are not allowed but unlockings of shared tasks can be done.
- ii) All the properties are same as that of rigorous 2PL, but it provides better concurrency.

DBMS NOTES BY RISHIRAJ

REFERENCE : KNOWLEDGE GATE YOUTUBE

Relational Algebra: After designing design of database, i.e. ER diagram design then converting it into relational model followed by normalization & indexing. Now, the task is ~~is~~ how to store, retrieve and modify data in database. Though here, we will be concentrating more on the ~~on~~ query part.

Query Language: Language in which the user requests some information from the database.

Procedural Query Language: Here, user instructs the system to perform a sequence of operations in order to produce ~~the~~ desired result. User tells what data to be retrieved & how to be retrieved.

Non-Procedural Query language: Here, user describes describes the desired information without giving the specific procedure for obtaining that information.

In practice, we use RDBMS (Practical implementation of ~~the~~ relational model)

- SQL is used to write query on it.
- Relational model is a conceptual/theoretical framework & RDBMS is its implementation.
- Relational algebra (procedural) & relational calculus (Non-procedural) are mathematical system or ~~the~~ query language used on Relational Model.

Fundamental Operators: ~~As~~ Relational algebra is one of the two formal query languages associated with relational algebra.

- Like any other mathematical system, it defines a no. of operators & use relations (tables) as operands.
- Every operator in relation algebra take one or two relations as input argument & generate a single relation as a result without a name.
- Relational Algebra does not consider duplicacy as it is based on set theory.
- In each query, we describe step by step procedural procedure for computing the desired output. So, it is a procedural query language.
- No use of English keywords.

Basic / Fundamental	Derived
Select ( $\sigma$ )	Natural Join ( $\sigma_{\alpha} \bowtie$ )
Project ( $\pi$ )	$\bowtie$ , $\bowtie[$ , $\bowtie]$
Union ( $\cup$ )	Intersection ( $\cap$ )
Set difference ( $-$ )	Division ( $\div$ )
<del>Cartesian Product (<math>\times</math>)</del>	
Rename ( $\delta$ )	

Select ( $\sigma$ ) : Select is a unary operator so that can take only one table at a time.

- It is a fundamental operator.
- Main idea of select is to find those tuples/rows in a relation which satisfies a given set condition.
- Syntax,  $\sigma$  condition/predicate (table-name)

- It has same function as of 'where' clause in SQL.
- Min. no. of tuples selected is 0.
- Max. no. of tuples selected is all the tuples.

Project ( $\Pi$ ): Project is a unary operator, takes one table at a time.

- Fundamental operator.
- Main idea of project is to select desired columns.
- Syntax,  $\Pi_{\text{column\_name}} (\text{table\_name})$
- It works as 'select' clause of SQL.

Operators  
Besides such as Union ( $\cup$ ), Intersection ( $\cap$ ) and Set Difference ( $-$ ) are needed for multiple table queries.

Cartesian Product ( $\times$ ): It is a binary operator, takes two tables at a time.

- Fundamental Operator.
- Allows us to combine information between two tables.
- $R_1(1, 2, 3, \dots, n), R_2(1, 2, 3, \dots, m) \rightarrow R_1 \times R_2(1, 2, 3, \dots, m+n)$

After performing cartesian product, we get redundant (orphus) tuples.

Q.) Find customer name having account balance < 100.

Sol:  $\Pi_{\text{cust\_name}} (G \underbrace{\text{Account.account\_no} = \text{Depositor.account\_no}, \text{balance} < 100}_{\text{This removes the redundancy}}) \text{ (Account } \times \text{ Depositor)}$

• In relational algebra, cartesian product is a commutative operator.

## Introduction to SQL

- Number of query languages are in use.
- Structured Query Language, SQL is the most popular & widely accepted query language.
- It is domain specific & not general purpose, used in design and management of data held in DBMS.
- Can do ~~query~~ much more than just a query on a database. Can define structure of database, modify data in database, specify security constraints & number of other tasks.
- Based on relational algebra and tuple relational calculus.

Parts of SQL :

i) Data Definition Language (DDL) - It provides commands for defining relational schemes, deleting relations and modifying relational schemas.

Integrity - Provides commands for integrity constraints constraints on the the data stored in database. e.g., primary key cannot be null, foreign key reference, etc.

View Definition - Commands for defining views, e.g. sorting the result according to some attribute.

Authorization - Commands for access rights, e.g., read only, R/W grants, etc.

ii) Data Manipulation Language (DML) - It provides the ability to control query information from the database and insert tuples, delete tuples and modify tuples in the database, e.g., insert, delete, etc.

iii) Transaction Control - SQL includes commands for specifying the begin beginning and ending of transactions,

S., commit, rollback, save points, etc.

### SQL Clauses

Select Query: i) For data retrieval in SQL, i/p & o/p both are relations.

- ii) No. of relation ip to query will be atleast one, but op will always be a single relation without any name unless specified.
- iii) Basic structure of SQL query consists of three clauses.
- iv) It should be noted that select and from are mandatory and if not required, then it is not necessary to write where.
- v) SQL in general, is ~~is~~ not case sensitive.
- vi) SQL allows duplicates in an ip relation as well as the in the result of SQL expression.

Select Clause: i) It is used to pick column required in result of the query, out of all the columns in the ip relation.

ii) Order in which columns are represented in the select clause will be same in the order of columns in the result.

iii) If all columns are required, then \* can be used.

iv) May also contain arithmetic expression involving operators like +, -, /, \*, etc. However, it does not change the data in the database.

\* DISTINCT keyword removes duplicacy.

Where Clause: i) It is used to specify conditions or predicates.

ii) It can have expressions involving comparison operators  $<$ ,  $>$ ,  $\leq$ ,  $\geq$ ,  $\neq$ , etc.

iii) SQL allows use of logical operators in where clause like AND, OR, NOT, etc.

iv) It allows 'between' comparison operator to specify where clause. It specifies  $\geq$  to some value and  $\leq$  to some other value. Similarly 'not between' is also present.

Set Operators : i) SQL provides set operations like Union, Intersect, except/minus/set difference which corresponds to the mathematical set operations  $\cup$ ,  $\cap$ , - respectively.

ii) Here, union, ~~intersection~~ ~~intersect~~, intersect, minus auto automatically eliminates duplicate tuples.

iii) If we want to retain duplicate tuples, we write 'ALL' after operators.

iv) Set operations are only applicable between two relations if they are comparable compatible - i.e. having same number of columns with correspondingly same domain.

Cartesian Product : i) It combines two relations into one by concatenation of each tuple of the first relation with every possible tuple of the second relation.

ii) From clause by itself defines the C.P. of the relation listed.

iii) C.P. is commutative in nature, so order doesn't matter.

iv) If same attribute name appears in both relations, we prefix the name of the relation from which which the attribute originally comes.

Natural Join : i) It will work same as cartesian product but consider only those pair of tuples with the same value of those attribute that appear in the schema of both relations.

- (ii) Notice that common attribute comes only once and <sup>prefix</sup> is also not required as we allow only some values.
- (iii) Commutative in nature.
- (iv) May lead to data loss if different attribute to be matched have different names, or some values occur only in one table.

Inner Join : It is an operation which works same as cartesian product if used alone.

- (i) But if join is used with keyword 'using', it provides additional functionality.
- (ii) Provides ability to explicitly choose columns which must be used by join for comparison & removal of redundant tuples.
- (iii) If there are more than one column common between two tables but we don't want each of them to be considered, then we specify which column to be considered.

Join with ON : (i) In SQL, join with 'ON' allows us to use a general predicate over the relations being joined.

- (ii) This predicate is written like a 'WHERE' predicate except for the use of the keyword 'ON' instead of 'WHERE'.
- (iii) Like C.P., we must specify ~~out~~ which attribute must be matched.
- (iv) More powerful as predicates of ~~as~~ where can also be written with 'ON'.
- (v) 'ON' condition can express any SQL predicate.
- (vi) It improves the readability of the query.
- (vii) In case of outer join, 'ON' condition behaves differently from 'where'.

Outer Join : One problem with natural join is that only those values

values that appear in both relations will manage to reach the final table. But if some data is explicitly in table one or table two, then it will be lost.

There are 3 forms of outer join :

- i) Left outer join (Left join)
- ii) Right outer join (Right join)
- iii) Full outer join (Full join / Complete join)

Rename Operation : i) SQL Aliases are used to give table or columns of a table, a temporary name.

- ii) Aliases only exist for the duration of the query.
- iii) Uses 'AS' alias clause & can appear both in 'select' and 'from' clauses.
- iv) Often used to name ; more readable, meaningful or smaller.
- v) Used for comparing a table with itself.
- vi) New temporary temporary name of table known as table alias / correlation variable / tuple variable.
- vii) Just creates a new name, but don't changes anything in the database.

String Operations :

Keep practicing SQL questions .

HackerRank and Leetcode