

Sequential files, index sequential files, direct files, Hashing, B-Trees, Index files, Relational models, Relational algebra and various operations, Relational and Tuple calculus.

X

X

- A database consists of a huge amount of data. The data is grouped within a table in RDBMS, and each table has related records. A user can see that the data is stored in form of data table, but in actual this huge amount of data is stored in physical memory in form of files.

### File:

A file is a named collection of related information that is recorded on secondary storage such as magnetic disks, optical disks and magnetic tapes.

### File organisation:

- File organisation refers to the representation or arrangement of data on external storage.
- or it refers to the logical relationships among various records that constitute the file. In simple terms, storing the files in certain order is called as file organisation.
- The methods by which records can be retrieved and updated from the files are called as access method.

## ↓ Types of file organisation

Engineering  
Notes  
File

### Sequential file organisation:-

The easiest method for file organisation is sequential file organisation method. In this method, the files are stored one after another in a sequential manner.

A sequential file contains records organised by the order in which they were entered. The order of records are fixed. It can only be read or written only sequentially.

(it occupies less space as records are stored in continuous order).

After you place a record into a sequential file; you cannot shorten, lengthen, or delete the record. However, you can update (REWRITE) a record if the length doesn't change. New records are added at the end of the file.

If the order in which you keep records in a file is not important then sequential file organisation is a good choice whether there are many records or only a few. Sequential output is also useful for printing reports.

→ There are two ways to implement this methods are as follows:-

### Pile File Method

It is quite simple and here, we store the records in sequence i.e. one after other in order in which they are inserted into the tables.

R <sub>1</sub>	R <sub>3</sub>	---	R <sub>5</sub>	R <sub>4</sub>
----------------	----------------	-----	----------------	----------------

↓  
Starting of  
file

↓  
end of  
file

### Sorted file Method

whenever a new record has to be inserted, it is always inserted in a sorted (ascending or descending) manner. It is based on primary or other key.

R <sub>1</sub>	R <sub>3</sub>	---	R <sub>7</sub>	R <sub>8</sub>	← R <sub>2</sub> new record
----------------	----------------	-----	----------------	----------------	-----------------------------

### advantages:

- fast and efficient methods for huge amount of data.
- simple design
- files are easily stored in magnetic tapes i.e. cheaper storage mechanism.

### disadvantages:

- Time wastage as we cannot jump on a particular record that is required, but we have to move in a sequential manner which takes our time.
- Sorting file method is inefficient as it takes time and space for sorting records.

### • Index sequential files:-

- Index sequential file organisation is an advanced sequential file organisation method.
- in this, files can be accessed sequentially as well as randomly with the help of the record key or (the primary key).
- it is possible to store the records in the middle of the file. It occupies more space.
- it also provides slow access (but it is fast as compared to sequential file organisation) as it takes time to search for the index.
- in this, records are written in sequential

Order but can be read only in sequential as well as random order.

- One or more keys can be created for storing and accessing the records.

→ The main drawback of sequential file is that searching operation is not efficient because in Sequential organisation, primary key of every record is compared with searching key. To optimise this operation concept of indexed sequential file is introduced.

- In this, a separate file for storing indexes of every record is maintained along with the master file.

eg:	IP	Position		Position	IP	name	Salary
	10	5		0			
	20	1		1	20	Alexa	2000
	30	7		2			
	40	3		3	40	Allu	4000
				4			
				5	10	Pam	1000
				6			
				7	30	Vay	3000
				8			

• no need to scan the entire memory block of record.

```

graph LR
    0 --> 10
    1 --> 20
    2 --> 30
    3 --> 40
    4 --> 10
    5 --> 20
    6 --> 30
    7 --> 40
  
```

(ISAM is a file management system developed at IBM that allows records to be accessed either sequentially or randomly).

### Advantages :-

- desired record can be accessed efficiently by using index which is maintained in separate file.
- Variable length records can also be handled using index sequential file.

### disadvantages :-

- atleast 2 files need to be maintained :-  
1) master file      2) Index file

- extra amount of memory is required in order to maintain index file.

### • Direct files :-

- direct access file is also known as random access organisation.

- in direct access file, all records are stored in direct access storage device (DASD), such as hard disk.

- The records are randomly placed throughout the file. The records does not need to be in sequence because they are updated directly and rewritten back in the same location.

- The file organisation is useful for immediate access to large amount of information.

→ application → OTP, railway reservation system, airline seat reservations.

- it is used in accessing large database. It is also known as hashing.

advantage :-

- It helps in online Transaction processing system (OLTP) like online railway reservation system.
- Sorting of the records is not required in this file organisation.
- It accesses desired records immediately, and can update several files quickly.
- It has better control over record allocation and it is more suitable for interactive online applications.

disadvantage :-

- direct access file does not provide backup facility.
- expensive hardware and software are required.
- File updating is more difficult when compared to that of sequential method.
- Less efficient use of storage space.
- Data may be accidentally erased or over written unless special precautions are taken.

- used for storing and retrieving data of a constant-time from database
- (also known as mapping technique of using hash function to map large volume to smaller values by using hashing)

Hashing :-

- Hashing is an efficient technique to directly search the location of desired data entries disk without using index structures.
- works well for large databases.

• Data is stored in the data blocks whose address is generated by using hash function. The memory location where these records are stored is called as data block or data bucket.

$$\text{Let } K = \text{key} \\ Q := K \bmod 10 \quad (K = 24) \\ \text{Q} = \text{Hash value} \\ \text{Memory for 10 Start from } 0, 1, 2, \dots, 9$$

① Data bucket -

data buckets are the memory locations where the records are stored. These buckets are also considered as unit of storage.  
Start the new record  
is used to generate address to a

② Hash function - (It is mathematical function which a hash function is a mapping function that maps all the set of given key to actual memory address. Generally, hash function has two primary key to generate the hash index - address of data block.)

③ Hash index -

The suffix of an entire hash value is taken as hash index. Every hash index has a depth value to signify how many bits are used for computing a hash function. The bits can address  $2^n$  buckets.

- used for store and retrieve data at a constant time from database
- (also known as mapping techniques as it tries to map large values to smaller values)

Hashing :- (by using hashing)

- Hashing is an efficient technique to directly search the location of desired data on the disk without using index structures.
- works well for large databases
- Data is stored in data blocks where address is generated by using hash function. The memory location where these records are stored is called as data block or data bucket.

$$\text{eq: } K \bmod 10 \quad (K=24)$$

↓ Hash file organisation

0
1
2
3

$$24 \bmod 10$$

④ = Hash value  
means 4 part 24 stored  
as 0000

① Data bucket -

data buckets are the memory locations where the records are stored. These buckets are also considered as unit of storage. Store the new record is used to generate data address to a

② Hash function - ( it is mathematical function which a

Hash function is a mapping function that maps all the set of search keys to actual record address. Generally, hash function uses the primary key to generate the hash index - address of data block.)

③ Hash index -

The prefix of an entire hash value is taken as hash index. Every hash index has a depth value to signify (how many bits are used for computing a hash function) The bits can address  $2^n$  buckets.

There are 2 types of Hashing:-

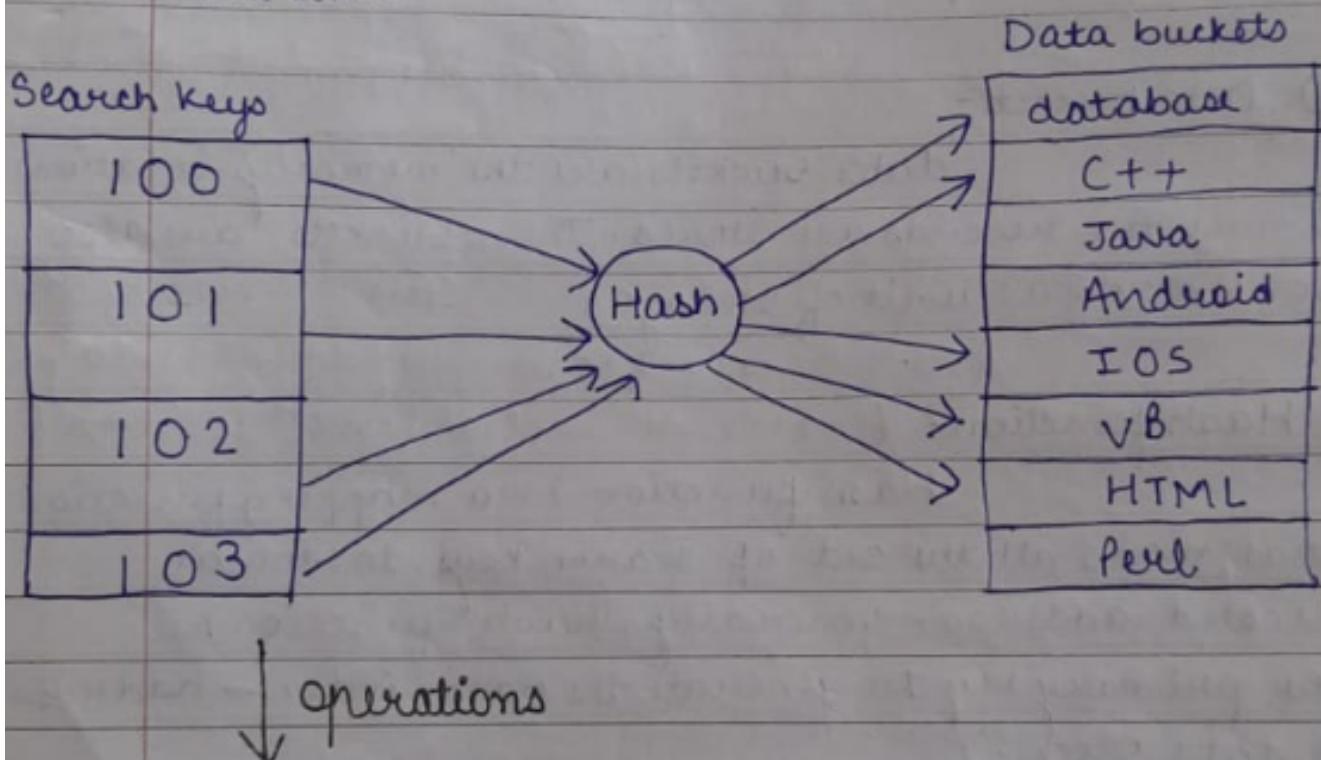
1) Static Hashing →

- in static hashing, when a search key value is provided, the hash function always computes the same address.

• for eg:-

if mod-4 hash function is used, then it shall generate only 5 values. The off addresses shall always be same for that function.

The number of buckets provided remains unchanged at all times.



① insertion: When a record is required to be entered using static hash, the hash function computes the bucket address for search key 'k', where the record will be stored.

$$\text{Bucket address} = h(k).$$

② delete: This is simply a search followed by a deletion operation.

⑤ Search: When a record needs to be retrieved, the same hash function can be used to retrieve the address of the bucket where the data is stored.

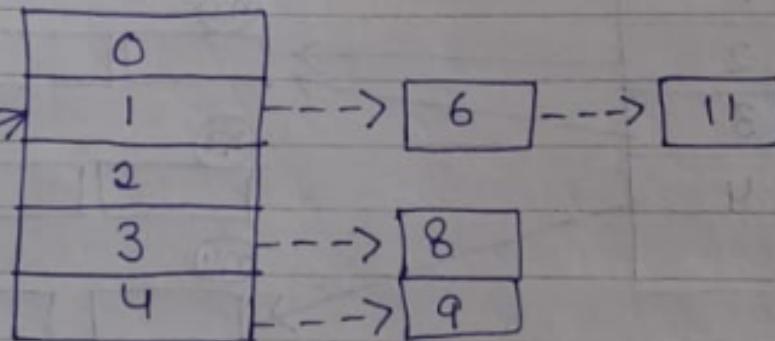
↓ Bucket overflow

The condition of bucket overflow is called as collision. This is a fatal state for any static hash function. In this case, overflow chaining can be used.

⑥ Overflow chaining:

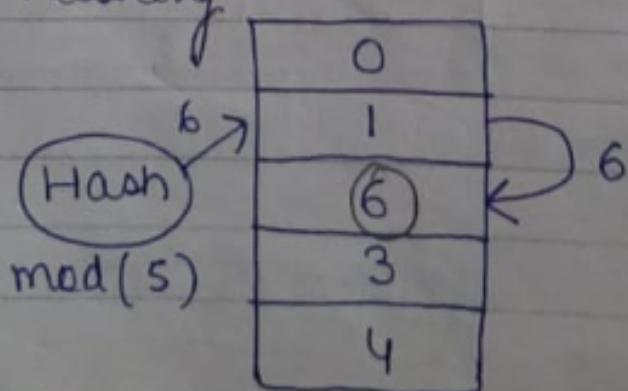
When buckets are full, a new bucket is allocated for the same hash result and is linked after the previous one. This mechanism is called closed hashing.

Data buckets



⑦ Linear Probing:

When a hash function generates an address at which data is already stored, the next free bucket is allocated to it. This mechanism is called open hashing.

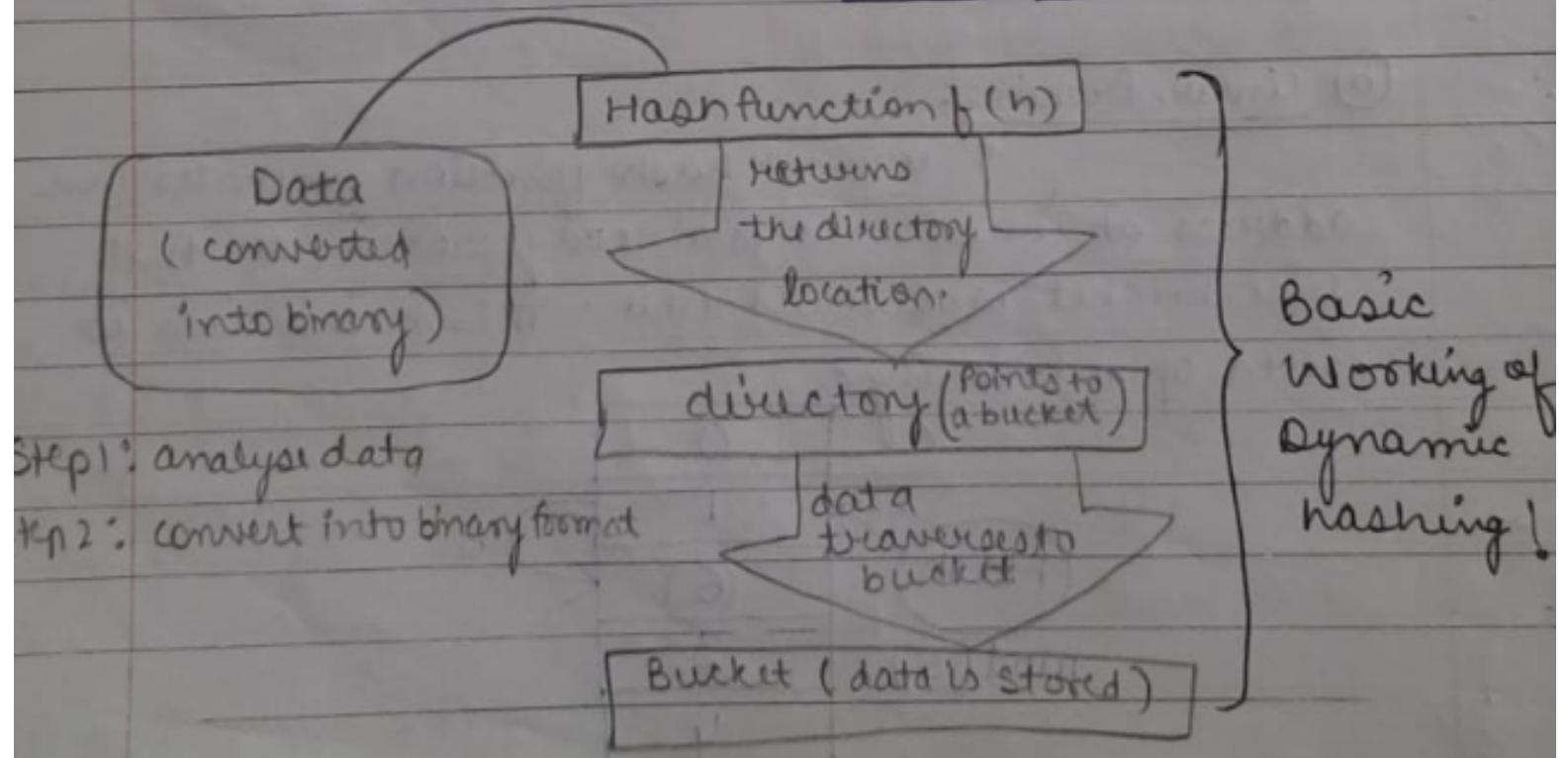
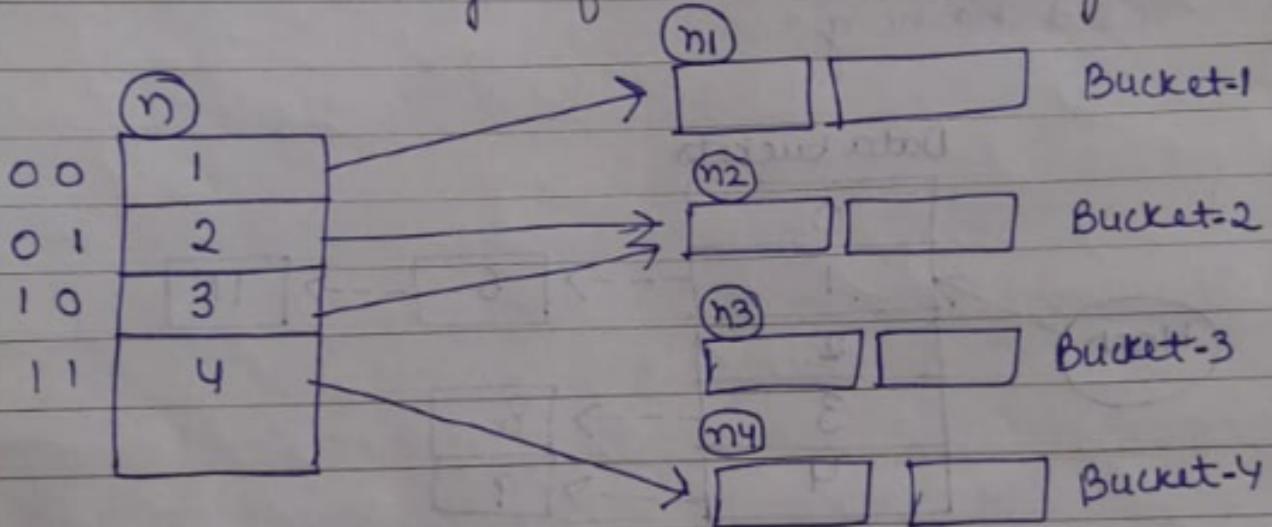


## 2) Dynamic hashing :

The problem with static hashing is that it does not expand or shrink dynamically as the size of the database grows or shrinks.

Dynamic Hashing provides a mechanism in which data buckets are added or removed dynamically and on-demand. Dynamic hashing is also known as extended hashing.

Hash function (in dynamic hashing), is made to produce a large number of values and only a few are used initially.



→ So, extended hashing is a dynamic hashing method where in directories and buckets are used to hash data. It is an aggressively flexible method in which the hash function also experiences dynamic change.

- Directories:

The directories store addresses of the bucket in pointer and it is assigned to each directory which may change each time when direct expansion takes place.

- Bucket:

The buckets are used to hash the actual data.

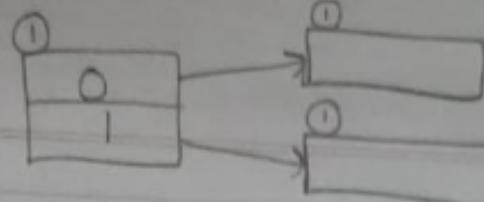
for eg → do extendable hashing for given data :-

16, 4, 6, 22, 24, 10, 31, 7, 9, 20, 26 .

↓ convert into binary

16 - 1 0000	7 - 00111
4 - 00100	9 - 01001
6 - 00110	20 - 10100
22 - 10110	26 - 11010
24 - 11000	
10 - 01010	
31 - 11111	

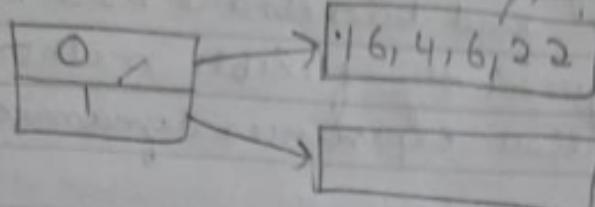
now assume the buckets :-



$$16 = 10000$$

insert 16, 4, 6, 22

only split  
0 as it is  
overflow



as order = 3,  
so it overflow!

$$4 = 00100$$

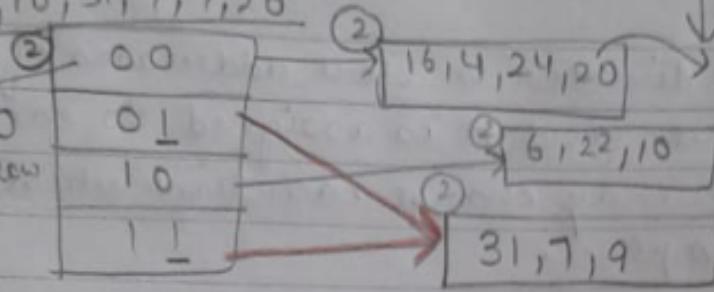
$$6 = 00110$$

$$22 = 10110$$

now,

input 24, 10, 31, 7, 9, 20

only split 00  
as it is overflow



$$24 = 11000$$

$$10 = 01010$$

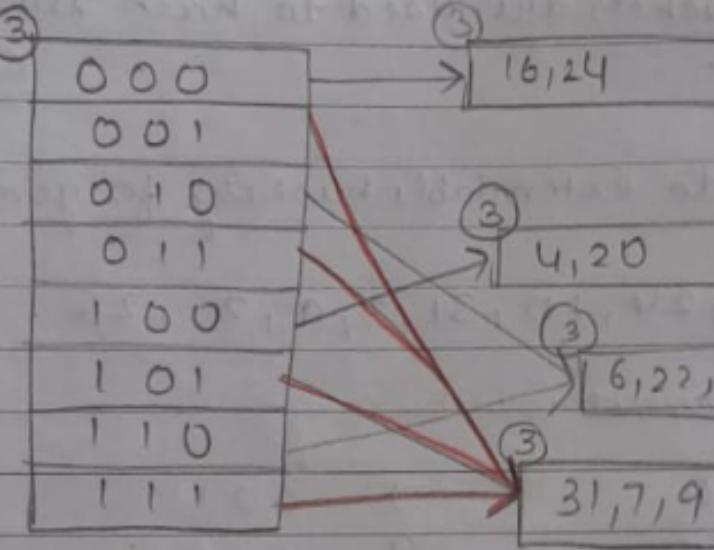
$$31 = 11111$$

$$7 = 00111$$

$$9 = 01001$$

$$20 = 10100$$

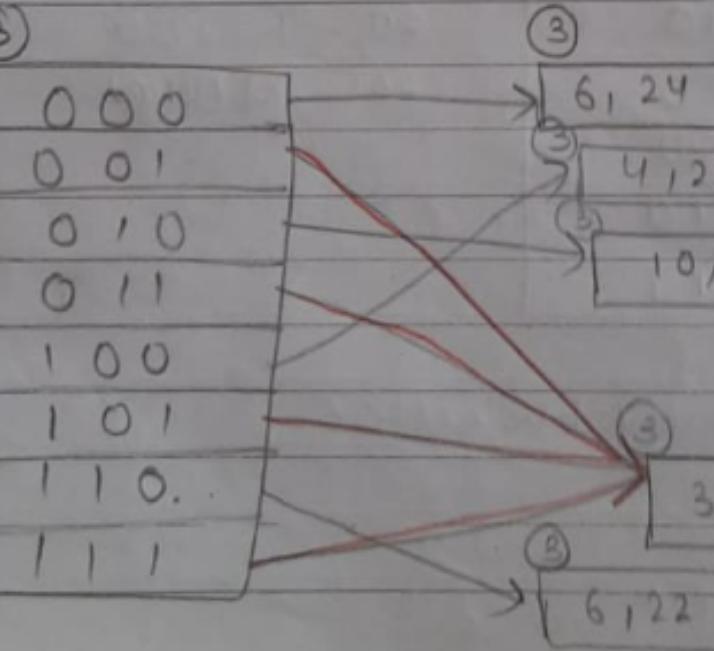
input 26



$$26 = 11010$$

overflow!

now,



$$6 = 00110$$

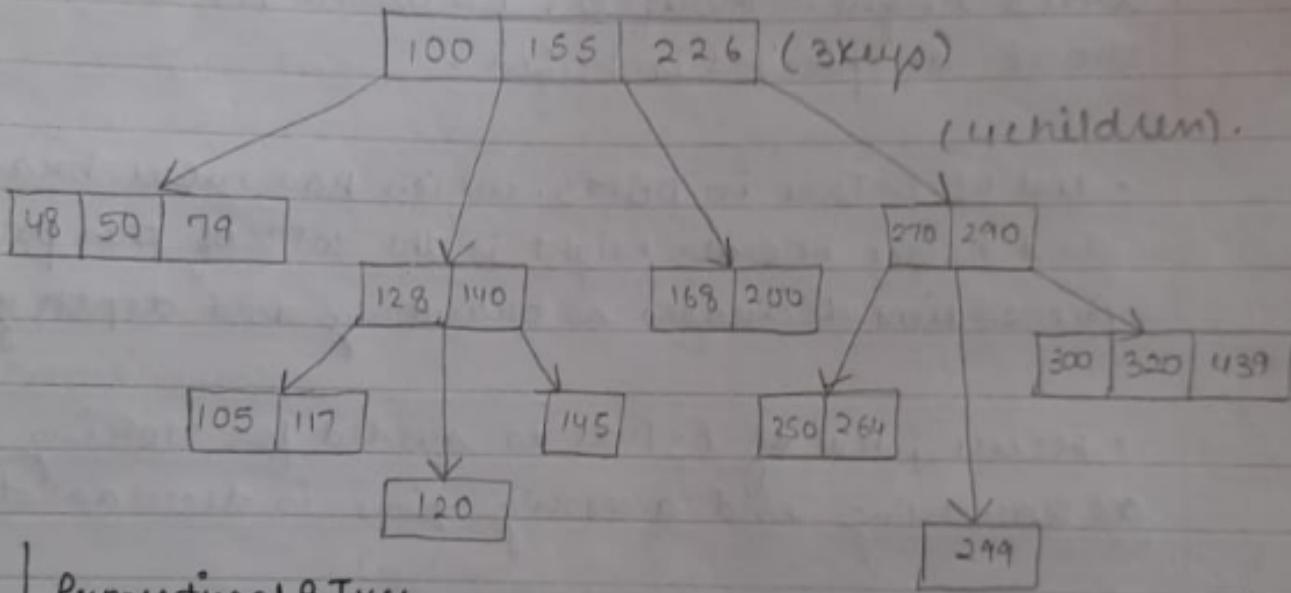
$$22 = 10110$$

$$10 = 01010$$

$$26 = 11010$$

extended  
washing!

- it has more than one key, & more than 2 children.
  - maintained sorted data.
  - all leaf nodes must be at same level.
  - every node has max.  $(m-1)$  keys.
- BTree index files:-**
- B Tree in DBMS is a  $m$ -way tree which self balances itself. due to their balanced structure, such trees are frequently used to manage and organize enormous databases and facilitate searches.
  - in B Tree, each node can have a maximum of  $n$  child nodes. eg → multilevel indexing.
  - Leaf nodes and indexing nodes will have both record references. B Tree is called as balanced stored trees as all the leaf nodes are at same levels.



### Properties of BTree

Following are some of the properties of B-Tree are as follows :-

- A non-leaf node's number of keys is one less than the number of its children.
- The number of keys in the root ranges from one to  $(m-1)$ .

Date / /

maximum. Therefore, root has a minimum of two and a max. of  $m$  children.

- The keys range from min  $\left[\frac{m}{2} - 1\right]$  to max.  $(m-1)$  for all nodes (non-leaf nodes) besides the root. Thus, they can have between  $(m)$  and  $\left(\frac{m}{2}\right)$  children.
- The level of each leaf node is the same.

↓ need of BTree

- For having optimized searching we cannot increase a tree's height. Therefore, we want the tree to be as short as possible in height.
- Use of B-Tree in DBMS, which has more branches and hence shorter height is the sol<sup>n</sup> of this problem. Access time decreases as branching and depth grow.
- Hence, use of B-Tree is needed for storing data as searching and accessing time is decreased.
- The cost of accessing the disc is high when searching tables. Therefore, minimising disc access is our goal.
- So to decrease time and cost, we use B-Tree for storing data as it makes the index fast.

## ↓ How database BTree indexing works?

- When BTree is used for database indexing, it becomes a little more complex because it has both a key and a value. The value serves as a reference to a particular data record. A payload is the collective term for the key and a value.
- For index data to particular key and value, the database first constructs a unique random index on a primary key for each of the supplied records. The keys and record byte streams are then all stored on a Btree. The random index that is generated is used for indexing of the data.
- So this indexing helps to decrease the searching time of data. In a B-Tree, all the data is stored on the leaf nodes; now for accessing a particular data index, database can make use of binary search on the leaf nodes as the data is stored in the sorted order.
- If indexing is not used, the database reads each and every records to locate the requested record and it increases time and cost for searching the records, so BTree indexing is very efficient.

## ↓ How searching happens in indexing database?

- The database does a search in the B-Tree for a given key and returns the index in  $O(\log(n))$  time. The record is then obtained by running a second B-tree search on

$O(\log(n))$  time using the discovered index. So, overall approx time taken for searching a record in a BTree in DBMS indexed database is  $O(\log(n))$ .

→ As in the previous example:-

The data is stored in sorted order according to the values, if we want to search for the node containing the value 48, so the following steps will be applied.

- First, the parent node with key having data 100 is checked, as 48 is less than 100 so the left children node of 100 is checked.
- In left children, there are 3 keys, so it will check from the leftmost key as the data is stored in sorted order.
- Leftmost element is having key value as 48 which matches the element to be searched, so that ~~how ever~~ the element we wanted to search is got.

↓ Conclusion :-

- A m-way tree that self-balances itself is called a "B-Tree".
- Due to their balanced structure, such trees are frequently used to manage and organize enormous

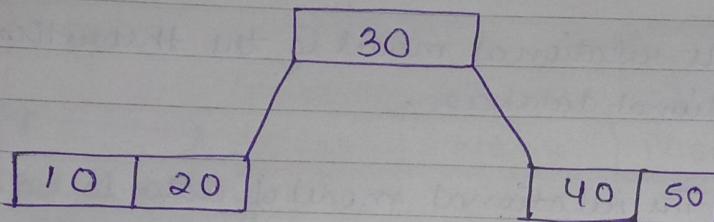
database and facilitate searches.

• B Tree is an example of multilevel indexing.

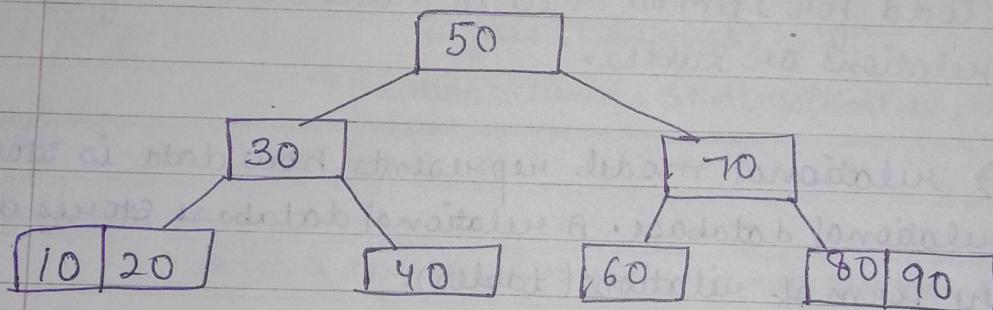
• uses of B Tree is needed for storing data as searching and accessing time is decreased.

• B Tree can be used for database indexing to improve searching and storing time of database.

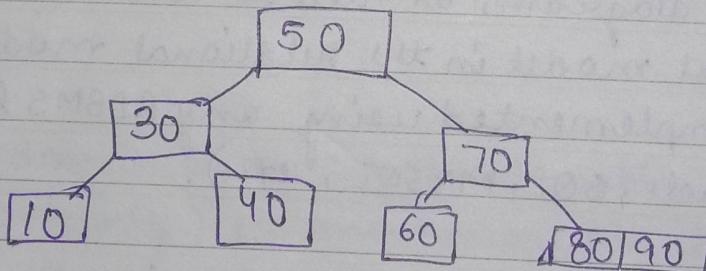
→ for eg: Construct B-Tree for search key values (10, 20, 30, 40, 50) when n=3.



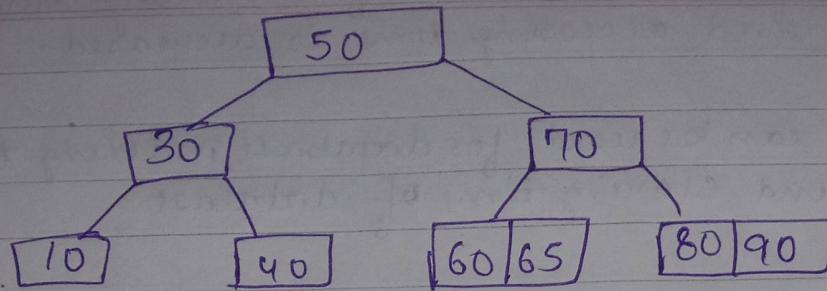
eg: values are 10, 20, 30, 40, 50, 60, 70, 80, 90.



a) Delete 20 from above tree :-



b) Insert 65 to above tree :-



- Relational Model:-

- The relational model is the theoretical basis of relational database.
- The relational model of data is based on the concept of relations.
- "Relational Model" is proposed by E-F. Codd for IBM in 1970 to model data in form of relations or tables.

→ relational model represents how data is stored in relational database. A relational database stores data in the form of relations (tables).

↓  
• after designing the conceptual model of database using ER diagram, we need to convert the conceptual model in the relational model which can be implemented using any RDBMS languages (like Oracle, SQL, MySQL, etc.).

• RDBMS stands for Relational database Management system. It is the basis for SQL, and for all model database systems like MS SQL server, IBM DB2, Oracle etc.

• So, RDBMS is a database management system (DBMS) that is based on relational model as introduced by E. F. Codd.

→ Relational model can be represented as table with columns (attribute or name) and rows (tuple).

eg:-

relation name			attribute / column
Primary key			
Tuple / row	Roll . NO	Name	Phone No.
	1	Ajay	9898373232
	2	Raj	9874444211
	3	Vijay	8923423411
	4	Aman	8886462644

Cardinality (no. of rows) = 4

Domain (integer : 0-9)

Degree (columns) = 3

Relation Schema = Students(Roll . NO, Name, Phone No.).

Terms

① relation:

it is a table with rows and columns.

② attribute:

an attribute is a named column of a relation.

③ Domain:

A domain is the set of allowable values (eg: Roll . NO can be 0-9 only) for one or more attributes.

④ Tuples:

A tuple is a row of a relation.

⑤ Relation Schema:

A relation schema represents the name of relation with its attributes (for eg: in this table, name of relation is student and attributes are roll.no., phone.no., name).

⑥ Relation instance:

A relation instance is a finite no. of set of tuples. It can never have duplicate tuples.

⑦ Degree:

The total number of columns or attributes in the relation.

⑧ Cardinality:

The total number of rows present in the table.

⑨ Relation Key:

Every row has one or multiple attributes, that can uniquely identify the row in relation (like roll.no.) which is called as Relation Key (primary key).

⑩ Tuple Variable:

It is the data stored in a record of the table.

- each relation has unique name; each row is unique; entries in column has the same domain; order of columns or rows is irrelevant; each cell of relation contains exactly one value.

Formal term	alternative 1	alternative 2
Relation	Table	File
Tuple	Row	Record
attribute	column	Field

- Relational algebra and various operations:-

- relational algebra is a procedural query language which takes a relation as an input and generates a relation as a output.
  - relational algebra is a language for expressing relational database queries.
  - It uses operator to perform queries.
- An operator used either be unary or binary.

→ types of operator are :- ① basic / fundamental operator, ② Additional / derived operator.

- Relational algebra operations work on one or more relations to define another relation without changing the original relations.
- Relational algebra operations are performed recursively on a relation.

→ for eg :-

I/P is a relation (table from which data has to be accessed) and O/P is also a relation (a temporary table holding the data asked for by the user).

we can  
use  
relational  
algebra  
to fetch data  
from this  
table (relation).

$\rightarrow I/P$

ID	name	Age
1	Akon	17
2	Bkon	19
3	Ckon	15
4	DKon	13

select name of  
students with age  
less than 17.

$\downarrow O/P$

Name
ckon
DKon

### various operations

#### ① Basic | Fundamental operations

- Selection ( $\sigma$ )
- Projection ( $\Pi$ )
- Union ( $\cup$ )
- Set difference ( $-$ )
- Cartesian product ( $\times$ )
- Rename ( $\rho$ )

#### ② Additional | Derived operations

- Natural Join ( $\bowtie$ )
- Set intersection ( $\cap$ )
- Division ( $\div$ )
- Assignment ( $\leftarrow$ )
- Left, right, Full outer Join  
( $\bowtie^L$ ,  $\bowtie^R$ ,  $\bowtie^F$ )

#### ① Selection :-

$\sigma$  (condition) it selects all tuples (row) that satisfy the selection condition from relation 'R'.

$$\sigma \text{ < condition } (R)$$

eg :-  $\sigma$  (age more than 17) (R)

② Projection :- Produce new relation with some of the attributes.

$$\pi_{\langle \text{attribute} \rangle}(R)$$

③ union :- produce a relation that includes all the tuples from R<sub>1</sub> or R<sub>2</sub>

$$(R_1 \cup R_2)$$

④ Set of difference :-  
 $(-) \quad$  Produces a relation that includes all tuples in R<sub>1</sub> that are not in R<sub>2</sub>?

$$R_1 - R_2$$

⑤ Cartesian product :-  
 also known as Cross product (X)  
 Produces A relation that has attributes of R<sub>1</sub> and R<sub>2</sub>, and includes as tuples, all possible combination of tuples from R<sub>1</sub> and R<sub>2</sub>.

$$R_1 \times R_2$$

⑥ Rename :-  
 $(P) \quad$  used to give a name to a relation obtained after applying any relational algebra operation

$$P(R, \epsilon)$$

⑦ Natural Join :-

$(\bowtie)$

Join two relation having any no. of attributes.

$R_1 \bowtie R_2$

⑧ Set intersection :-

$(\cap)$

produces a relation that includes all the tuples in both  $R_1$  &  $R_2$ .

$R_1 \cap R_2$

⑨ Division :-

$(\div)$

queries that includes 'FOR ALL', 'EVERY'.

$R_1 \div R_2$

→ for eg:-

### Selection

Student

Roll No.	Name	Age	Address
1	A	20	Bhopal
2	B	17	Mumbai
3	C	16	Mumbai
4	D	19	Delhi
5	E	18	Delhi

① Select student whose name is D

$\sigma \text{ name} = "D" (\text{student})$

Roll No.	Name	Age	Address
4	D	19	Delhi

② Select students whose age is greater than 17

$\sigma \text{ age} > 17 (\text{student})$

Roll No.	Name	Age	Address
1	A	20	Bhopal
4	D	19	Delhi
5	E	18	Delhi

### Projection

① Display (or project) name of students in student table.

Name
A
B
C
D
E

$\Pi_{<\text{name}>}(\text{student})$

② display address of students in student table.

address
Shopal
Mumbai
Delhi

Note:

→ by default, projection removes duplicate values.

Surf

eg :-

student	
Roll No.	Name
1	A
2	B
3	C
4	D

employee	
emp. no.	Name
2	B
8	G
9	H

union ①  $\Pi_{<\text{name}>}(\text{student}) \cup$

②  $\Pi_{<\text{student}>} \cup \Pi_{<\text{employee}>}$

$\Pi_{<\text{name}>}(\text{employee})$

Roll No.	Name
1	A
2	B
3	C
4	D
8	G
9	H

Name
A
B
C
D
G
H

### Set of difference

① (student) - (employee)

②  $\Pi_{\text{Name}}(\text{student}) - \Pi_{\text{Name}}(\text{employee})$

Name
A
C
D

Roll no.	name
1	A
3	C
4	D

### Cartesian product

①  $R_1 \times R_2$

e.g.:  $R_1$        $R_2$

	A	B	C	D	E
$\alpha$	1	$\alpha$	10	a	
$\beta$	2	$\beta$	10	a	
		$\beta$	20	b	
2		$\gamma$	10	b	
			4		

A	B	C	D	E
$\alpha$	1	$\alpha$	10	a
$\alpha$	1	$\beta$	10	a
$\alpha$	1	$\beta$	20	b
$\alpha$	1	$\gamma$	10	b
$\beta$	2	$\alpha$	10	a
$\beta$	2	$\beta$	10	a
$\beta$	2	$\beta$	20	b
$\beta$	2	$\gamma$	10	b

②  $\sigma_{A=c}(R_1 \times R_2)$

A	B	C	D	E
$\alpha$	1	$\alpha$	10	a
$\beta$	2	$\beta$	10	a
$\beta$	2	$\beta$	20	b

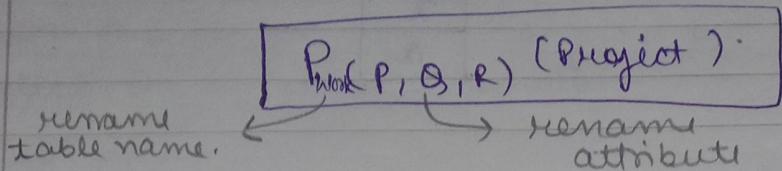
$$2 \times 4 = 8$$

### Renaming

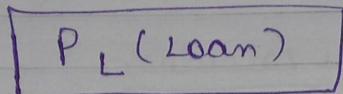
① query to rename attributes  
Name, Age of table Person to N, A,

$\rightarrow P_{(N, A)}(\text{Person})$

② query to rename the table name Project to work and its attribute to P,Q,R.



③ query to rename table name Loan to L.



Set intersection

① (student)  $\cap$  (employee)

Roll-NO.	Name
2	B

②  $\Pi_{\text{name}} (\text{student}) \cap \Pi_{\text{name}} (\text{employee})$

Name
B

Natural Join

① course  $\bowtie$  HOD

Courses			HOD	
CID	Course	Dept	Dept	Head
CS01	Database	CS	CS	A
ME01	Mechanics	ME	ME	B
EE01	electronics	EE	EE	C

CID	Course	Dept	Head
CS01	Database	CS	A
ME01	Mechanics	ME	B
EE01	electronics	EE	C

division ①  $A \div B$

A	B	$\frac{x}{a}$
x 4	y	
a 1	1	
b 2	2	
a 2		
d 4		

- Relational and Tuple calculus:-

→ Relational calculus →

• it is a non-procedural query language (or declarative language). It uses mathematical predicate calculus (or first order logic) instead of algebra.

• It tells what to do but never explains how to do.

• It provides description about the query to get the result whereas a relational algebra gives the method to get results.

→ when applied in database, it comes in two flavors:-

① Tuples Relational Calculus

• proposed by Codd in 1972

• works on tuples (or rows).

like SQL (structured query language).

② Domain Relational Calculus

• proposed by Lacroix and Pirotte in 1977

• works on domain of attributes (or columns).

like QBE (query by example)

- Tuple relational calculus:-

- A tuple relational calculus is a non-procedural query language that specifies to select the tuples in a relation.

- It can select the tuples with a range of values or tuples for certain attribute values etc. The resulting relation can have one or more tuples.

- TRC is variables range over tuples like SQL.

→ This query cannot be expressed using membership cond<sup>n</sup>.

- It is used for selecting the tuples in a relation that satisfies the given cond<sup>n</sup> (or predicate).

notation :-  $\{ T \mid P(T) \}$  or  $\{ T \mid \text{condition}(T) \}$



where T is resulting tuples and P(T) is a cond<sup>n</sup> used to fetch T.

for eg:-

$\{ T \mid \text{employee}(T) \text{ and } T.\text{Dept-ID} = 10 \}$

This selects all the tuples of employee names whose work for department 10.

- Domain relational calculus:-

- A domain relational calculus uses the list of attributes to be selected from the relation based on the condition.

- it is the same as TRC but differs by selecting the attributes rather than selecting whole tuples.
- DRC is a variable range over domain elements like query - by - example (QBE).

notation :-

$$\{ a_1, a_2, \dots, a_n | P(a_1, a_2, \dots, a_n) \}$$



where  $a_1, a_2, \dots, a_n$  are attributes of the relation and  $P$  is the cond<sup>n</sup>.

for eg :-

$$\{ | < \text{EMPLOYEE} > \text{Dept-ID} = 10 \}$$

Select EMP\_ID and EMP\_NAME of employees who work for department 10.

The query can be expressed using a membership condition.

X

X