

## UNIT-II

### Dictionary

Def:- Dictionary is a collection of ~~pair~~ of key value pairs. Every key value is associated with the corresponding key.

#### Advantages of dictionaries:

- ① Speed : By using dictionaries increase the process speed for some type of data.
- ② Memory efficient :- By using dictionaries we can efficiently store ~~for~~ the data for trees and hash tables.
- ③ flexibility - Easily extend and modify the existing dictionaries.
- ④ unordered - ~~\*~~

#### Disadvantages of dictionaries:-

- ① Complexity :- To implement dictionary little bit complex.
- ② Space Complexity :- Depending on the size of the dictionary and chosen data structure, the memory requirement can be quite high.

Example for dictionary.

following example key is the roll number and his marks score as value.

Key value pair is {rollno; Score}  
key value.

{545: 80, 580: 85, 510: 80}

In key value pair key must be unique.

If key is already existed it can be update with new value.

If key is not existed it can be added as new dictionary element.

### Dictionary representation:-

There are two types of representations

① Linear list representation.

② Skip list representation.

### Linear List representation:-

It is almost same as linked list, here every node maintain the key value pair.

Key	value
Rollno 510	attn 10

510	10	NULL
-----	----	------

struct node

{

int key;

int value;

struct node \*next;

typedef struct node dnode;

}

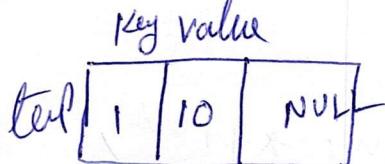
	insertion	removal	retrieval	②
unsorted array	$O(1)$	$O(n)$	$O(n)$	o
unsorted list	$O(1)$	$O(n)$	$O(n)$	
sorted array	$O(n)$	$O(n)$	$O(\log n)$	
sorted list	$O(n)$	$O(n)$	$O(n)$ ,	

## Operations on dictionaries:-

- ① Insert / update
- ② Delete
- ③ Retrieve / search
- ④ Traverse / display.

Insert :- Consider initially the dictionary is empty,  
that means start is NULL.

Create a new node with some key value.

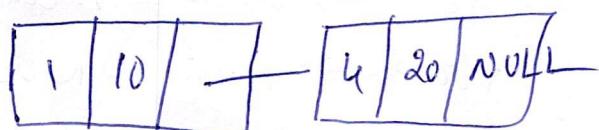


start

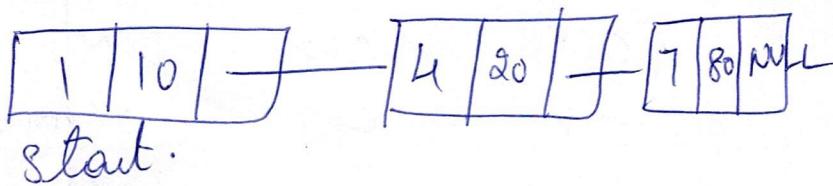
now temp node become start.

Q. If dictionary is not empty find appropriate position then insert.

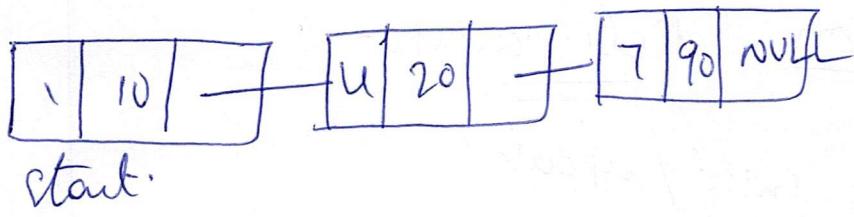
now insert (4, 20)



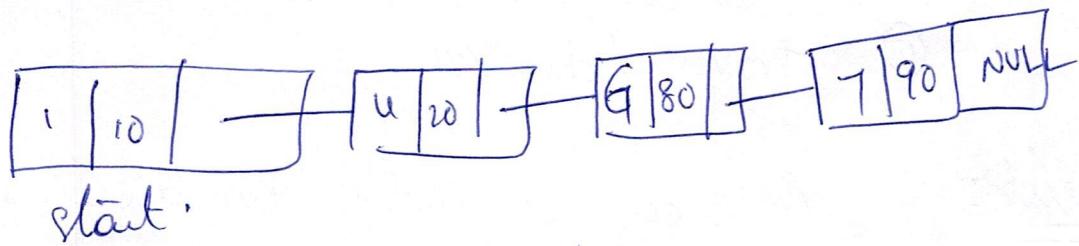
now insert (7, 80)



now insert (7, 90)



now insert (6, 80)



algorithm insert (key, value)

{  
// start is global identifier represents the  
head node. ~~start~~

dnode \*temp, \*curr, \*prev;

curr = prev = start;  
if ( start == NULL )

{

start = malloc( sizeof(dnode) );  
start → key = key;  
start → value = value;  
start → next = NULL;

else {

⑧

$cur = prev = start$

while(( $cur \rightarrow key != key$ ) && ( $cur \rightarrow next != NULL$ ) &&  
 {  
 ( $cur \rightarrow key < key$ ))}

$prev = cur;$

$cur = cur \rightarrow next;$

{}

// update with existing key

if ( $cur \rightarrow key == key$ )

$cur \rightarrow value = value;$

else

if ( $cur \rightarrow key > key$ ) // insert at middle position.

{  
~~create~~ temp = malloc(sizeof(node));  
~~temp~~ temp  $\rightarrow$  key = key; temp  $\rightarrow$  val = value;

$prev \rightarrow next = temp;$

$temp \rightarrow next = cur;$

}

else

if ( $cur == start \& cur \rightarrow key > key$ )

{

// create node temp;

$temp \rightarrow next = start;$

$start = temp;$

}

else

{

$cur \rightarrow next = temp;$

$temp \rightarrow next = NULL;$

}

5

3

delete (key)

{

// start is global Variable

cur = prev = start;

if ( cur == NULL )

    Put (" dictionary is empty");

else

{

    while ( cur->key != key && cur->next != NULL )

{

        prev = cur;

        cur = cur->next;

}

    if ( cur == start )

{

        start = start->next;

        free ( cur );

}

else

    if ( cur->next == NULL && cur->key != key )

        Put ("element not found");

else

    if ( cur->key == key ) && ( cur->next == NULL )

{

        prev->next = NULL;

        free ( prev );

}

{

        prev->next = cur->next;

        free ( cur );

3

3

3

## Traverse/display()

{  
// start is global identifier to represent  
start node of the dictionary.

if ( start == NULL )

{

    Put ("Dictionary is empty");

else {  
    3

    start temp = start;

    while ( temp != NULL ).

{

        Put ("temp->key , temp->value");

        3 temp = temp->next ;

}

5

## search/retrieve (key)

{

// start --

if ( start == NULL )

    Put ("dictionary is empty");

else {

{

    curr = prev = start;

    while ( curr->next != NULL && curr->key != key )

{

        prev = curr;

        3 curr = ~~curr~~ curr->next ;

~~if ( cur → key == key )~~

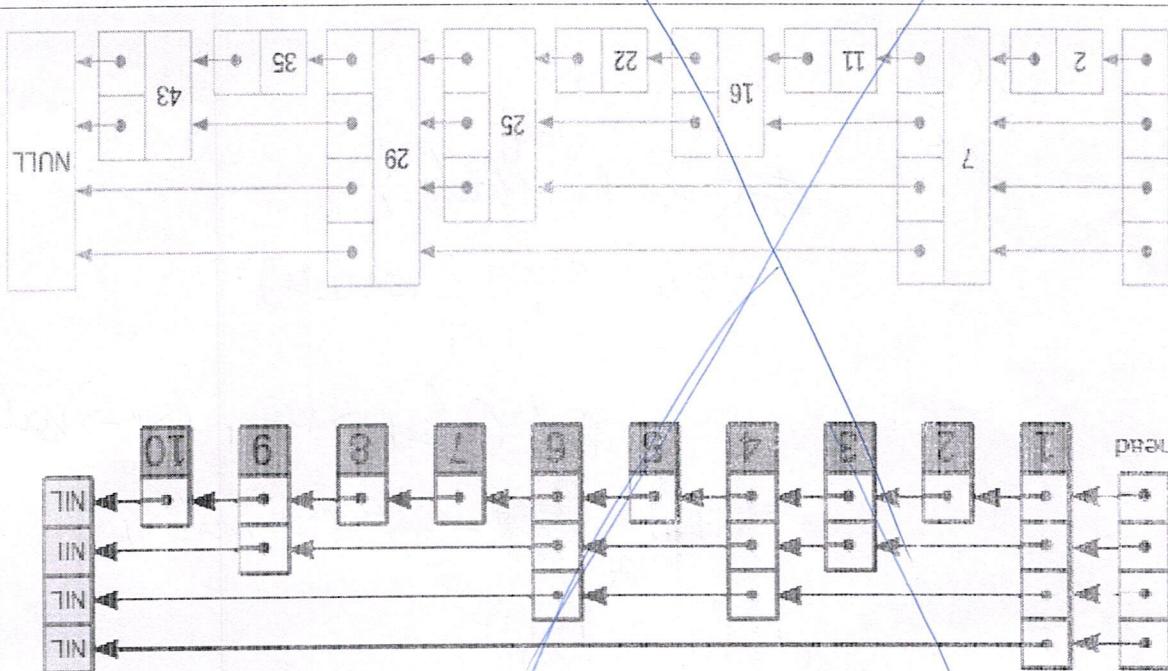
~~cout << "key found and its value is " << cur → val;~~

~~else~~

~~cout << "key not found";~~

A skip list is built up of layers. The lowest layer (i.e. bottom layer) is an ordinary ordered linked list. The higher layers are like 'express lane' where the nodes are skipped (observe the above figure).

### Structure of Skip List



A skip list allows the process of item look up in efficient manner. The skip list data structure skips over many of the items of the full list in one step, that's why it is known as skip list.

A skip list is a data structure that is used for storing a sorted list of items with a help of hierarchy of linked lists that connect increasingly sparse subsequences of the items.

### Skip List

if ( $cur \rightarrow key = key$ )

    Put ("key found and its value is ",  $cur \rightarrow value$ );

else

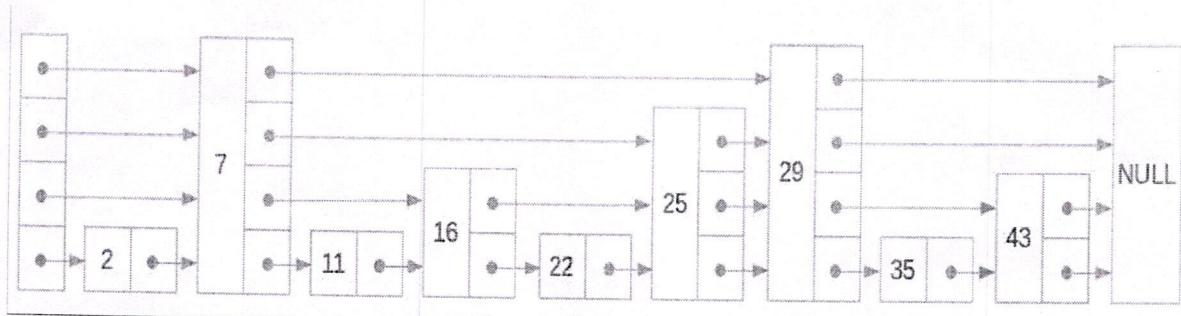
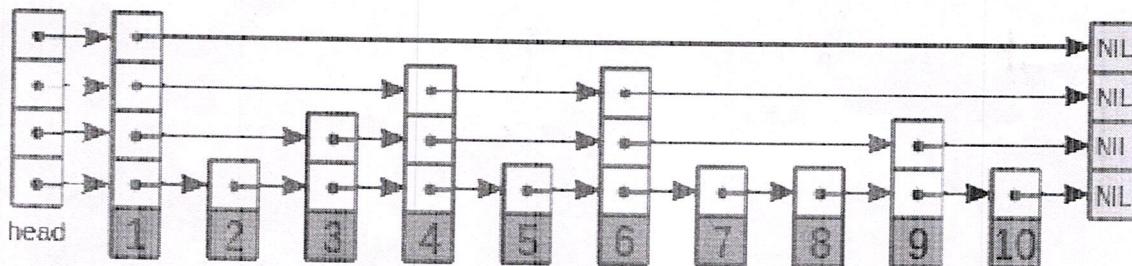
    Put ("key not found");

⑤

## Skip List

A skip list is a data structure that is used for storing a sorted list of items with a help of hierarchy of linked lists that connect increasingly sparse sub sequences of the items.

A skip list allows the process of item look up in efficient manner. The skip list data structure skips over many of the items of the full list in one step, that's why it is known as skip list.



### Structure of Skip List

A skip list is built up of layers. The lowest layer (i.e. bottom layer) is an ordinary ordered linked list. The higher layers are like 'express lane' where the nodes are skipped (observe the above figure).

## **Advantages of the Skip list**

1. The skip list is simple to implement as compared to the hash table and the binary search tree.
2. It is very simple to find a node in the list because it stores the nodes in sorted form.
3. The skip list is a robust and reliable list.

## **Disadvantages of the Skip list**

1. It requires more memory than the balanced tree.
2. Reverse searching is not allowed.
3. The skip list searches the node much slower than the linked list.

## **Applications of the Skip list**

1. It is used in distributed applications.
2. It is also used with the QMap (Open Multimedia Application Platform) template class.
3. The indexing of the skip list is used in running median problems.

## **Searching Process**

When an element is tried to search, the search begins at the head element of the top list. It proceeds horizontally until the current element is greater than or equal to the target. If current element and target are matched, it means they are equal and search gets finished.

## **Skip List Basic Operations**

There are the following types of operations in the skip list.

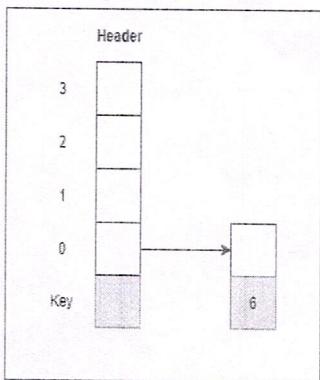
- **Insertion operation:** It is used to add a new node to a particular location in a specific situation.
- **Deletion operation:** It is used to delete a node in a specific situation.
- **Search Operation:** The search operation is used to search a particular node in a skip list.

### Insertion operation:

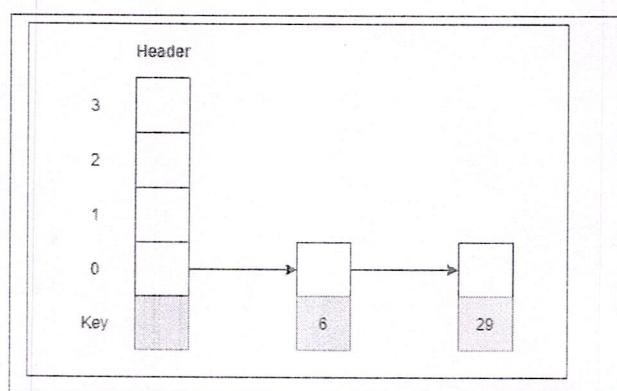
**Example 1:** Create a skip list, we want to insert these following keys in the empty skip list.

- 1. 6 with level 1.
- 2. 9 with level 3.
- 3. 29 with level 1.
- 4. 17 with level 1.
- 5. 22 with level 4.
- 6. 4 with level 2.

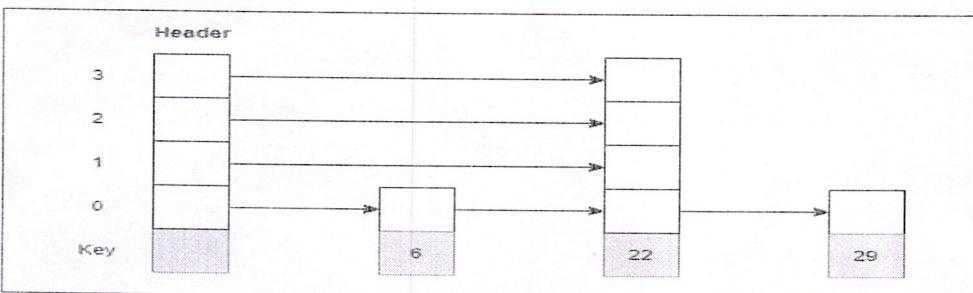
6 with level 1.



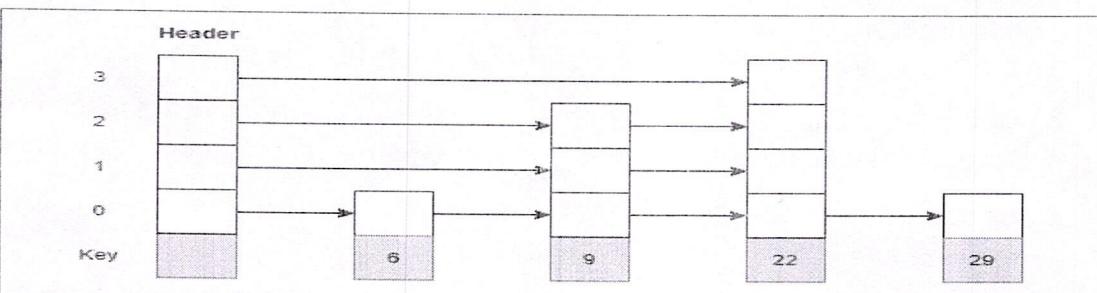
29 with level 1.



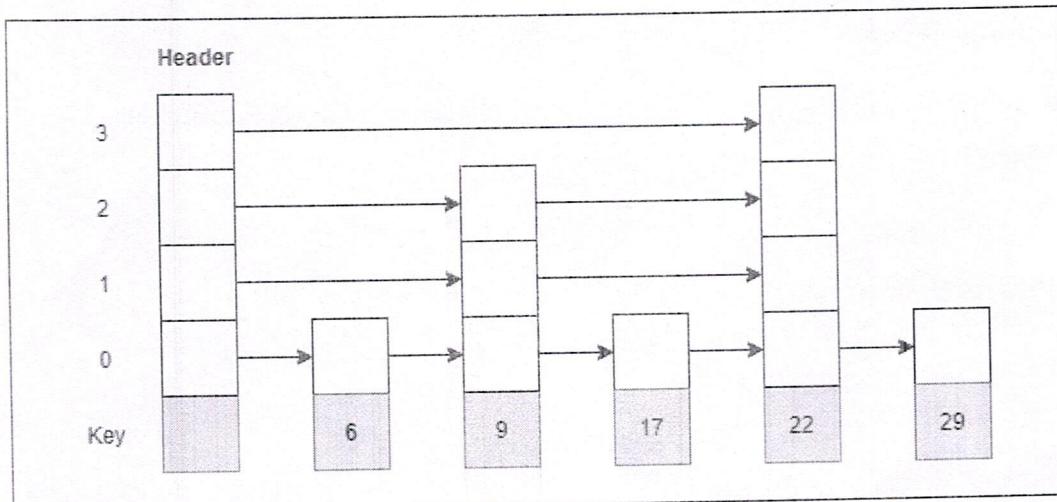
Insert 22 with level 4



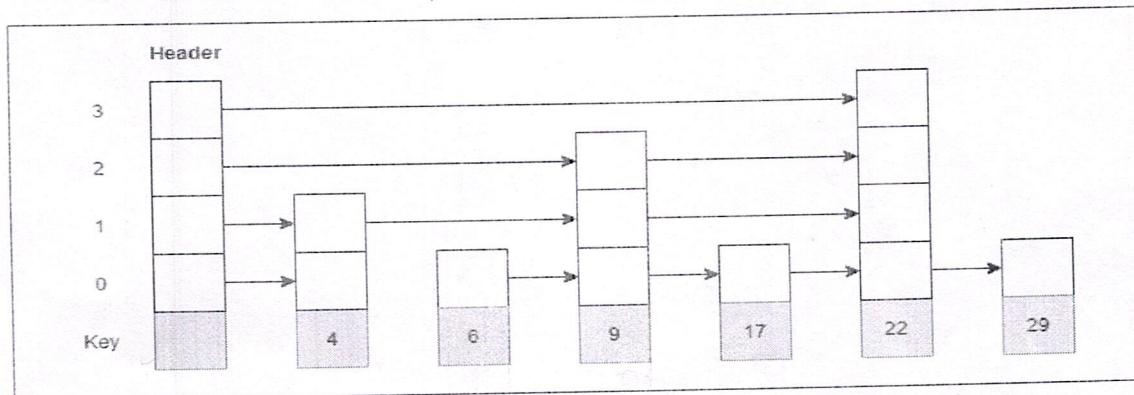
9 with level 3.



Insert 17 with level 1



Insert 4 with level 2



### Algorithm steps to insert element into skip list.

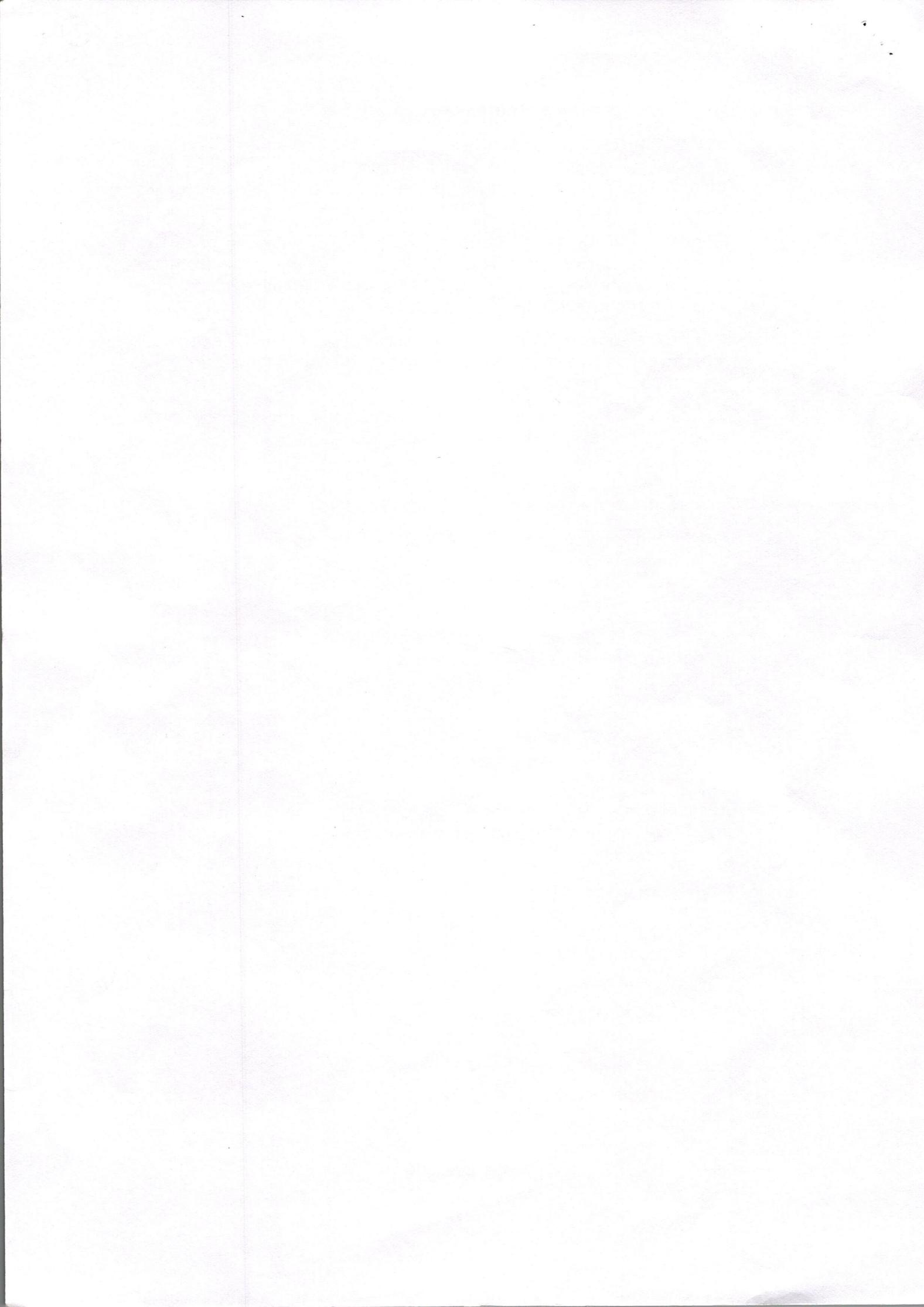
1. Create a new node with the value you want to insert.
2. Determine the level of the new node. You can do this randomly by flipping coins or using some other probabilistic method. The higher the node's level, the fewer nodes it will skip, leading to faster access times. Typically, the maximum level of a skip list is around  $\log_2(n)$ , where n is the number of elements in the list.
3. Start at the top-left corner (head) of the skip list and traverse the list horizontally while moving to the right until you find a node with a smaller value or reach the end of the list at the current level.
4. When you find the right spot, insert the new node into the current level's linked list by updating the pointers of the adjacent nodes.
5. Repeat this process for each level, going down one level at a time until you reach the base level.
6. The insertion is complete, and the new node is successfully added to the skip list

**Algorithm steps to delete element from skip list.**

1. Start at the top-left corner (head) of the skip list and traverse the list horizontally while moving to the right until you find the node with the target value or reach the end of the list at the current level.
2. When you find the node with the target value, you need to update the pointers of the adjacent nodes to remove the node from the skip list. To do this, you'll need to keep track of the update nodes (similar to the insertion algorithm).
3. Once you've updated the pointers at the current level, repeat this process for each level, going down one level at a time, until you reach the base level.
4. After completing the deletion for all levels, the node with the target value is successfully removed from the skip list.

**Algorithm steps to delete element from skip list.**

1. Start at the top-left corner (head) of the skip list and traverse the list horizontally while moving to the right until you find the node with the target value or reach the end of the list at the current level.
2. When you find the node with the target value, you need to update the pointers of the adjacent nodes to remove the node from the skip list. To do this, you'll need to keep track of the update nodes (similar to the insertion algorithm).
3. Once you've updated the pointers at the current level, repeat this process for each level, going down one level at a time, until you reach the base level.
4. After completing the deletion for all levels, the node with the target value is successfully removed from the skip list.



## Hashing

(8)

Hash :- Hashing is technique to map certain data to a particular key. In a large database, the data can be quickly stored and retrieving using the key.

Hash Table :- We store the identifiers in a fixed size table called hash table.. we use the arithmetic function  $f$  to determine the address or location.  
So  $f(x)$  gives the hash or home address of  $x$  in the table.

## Hash functions:-

A hash function  $f$  transforming an identifier  $x$  into a bucket address in the hash table. we need a hash function that is easy to compute and that minimizes the number of collisions.

Type of hash functions! - There are different types of hash functions to find/calculate the hash bucket address.

- ① Division method
- ② Multiplication method
- ③ Mid Square method
- ④ Digit folding method
- ⑤ Digit analysis method

Division method :- This method used to modulo operator to compute the hash code for a key. The modulo operator returns the remainder of the key.

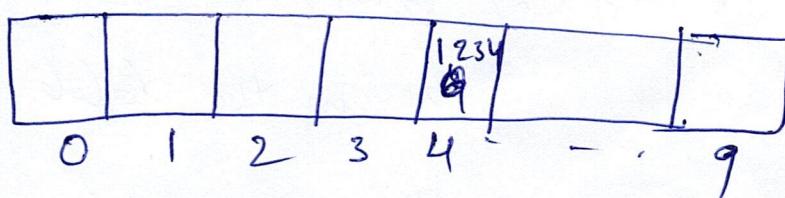
Example : The size of the array is - 10  
and key is 1234

division method function

$$f(\cancel{\text{key}}) = \text{key} \% . \text{size}$$

$$f(1234) = 1234 \% 10$$

$$= \underline{\underline{4}}$$



Note :- size of the array take it as prime number so it can reduce the number of collisions.

Advantage :- It gives the fast access.

multiplication method :-

This method involves the key by a constant value then taking the fractional part of the result. The fractional part is multiplied by the size of the table, and then take floor of that value taken as the hash code.

$$h(\text{key}) = \text{floor}(\text{size}(\text{key} \bmod 1))$$

(9)

K is the key

A is constant - 0.6180339887

size is size of the table. - 13

$$\begin{aligned} h(123) &= \text{floor}(13((123 + 0.6180339887 \% 1) \\ &\quad \xrightarrow{\text{floor}(76.01818)}) \\ &= \text{floor}(13 * (0.01818)) \\ &= \text{floor}(\cancel{13} * 0.236848) \\ &= \underline{\underline{0}} \end{aligned}$$

$$\begin{aligned} h(125) &= \text{floor}(13((125 + 0.6180339887 \% 1))) \\ &= \text{floor}(13(77.25425 \% 1)) \\ &= \text{floor}(13(0.254249)) \\ &= \text{floor}(3.305232) \\ &= \underline{\underline{3}} \end{aligned}$$

- Advantages:
- ① It produces the good distribution of keys.
  - ② It is suitable for both real and integer keys.

- Disadvantages:
- ① It is more complex to implement.
  - ② It may not work for small table size.
  - ③

## mid square method :-

This method involves squaring the key and extracting some digits from middle of the generated value.

$$h(\text{key}) \rightarrow \text{mid}(K * K);$$

$$h(123) = \text{mid}(23 * 123)$$

$$= \text{mid}(15129)$$

$$= \underline{\underline{1}}$$

$$h(56) = \text{mid}(56 * 56)$$

$$= \text{mid}(3136)$$

$$= \underline{\underline{1}}$$

$$h(158) = \text{mid}(158 * 158)$$

$$= \text{mid}(24964)$$

$$= \underline{\underline{9}}$$

folding method :- The key value  $K$  should be divided into a specific number of parts, such as  $K_1, K_2, \dots, K_n$ . Add each component separately.

$$K = K_1, K_2, \dots, K_n.$$

$$h(K) = K_1 + K_2 + K_3 + \dots + K_n.$$

$$h(K) = S.$$

Example :-

$$K = 1054321$$

$$K_1 = 54 \quad K_2 = 32 \quad K_3 = 1$$

$$S = K_1 + K_2 + K_3$$

$$= 54 + 32 + 1$$

$$= 87$$

$$h(K) = 87.$$

Advantage :- independent distribution in hash table

disadvantage :-

It produces more collisions.

### Collision resolution Techniques :-

What is collision :- Since hash function gets us a small number of a key which is a big integer or string, there is possibility that two keys result in the same value. The situation where a newly inserted key maps to an already occupied slot in the hash table is called collision.

To handle this situation by using collision resolution methods.

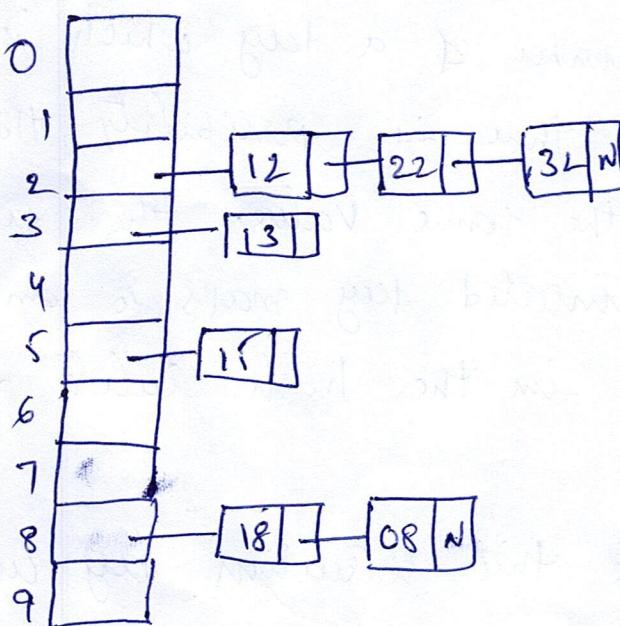
There are different types of collision resolution techniques

- ① Separate chaining
- ② Open addressing
  - a) linear probing
  - b) quadratic probing
  - c) double hashing
- ③ Rehashing
- ④ Extendible hashing.

### Separate Chaining:-

Separate Chaining is the collision resolution technique that is implemented using linked list. When two or more elements are hash to the same location, these elements are represented into a single linked list, like a chain. It is also known as open hashing.

12, 13, 15, 18, 08, 22, 32



$$h(\text{key}) = \text{key \% 10}$$

Advantages:-

- ① It is one of the simplest method to implement.
- ② We can add any number of elements to the Chain.

- Disadvantages:-
- ① The keys in the hash table are unevenly distributed.
  - ② Some amount of space wastage.
  - ③ Complexity in searching is  $O(n)$  in worst case.

### Open Addressing:- (closed hashing)

In open addressing, all elements are stored in the hash table itself. so at any point, the size of the table must be greater than or equal to the total number of keys. This approach is also known as closed hashing.

There are different ways to implement open addressing.

- ① Linear probing.
- ② Quadratic probing
- ③ Double hashing.

Linear Probing:- Linear probing collisions are resolved by searching the hash table consequently until empty cell is found. This process is called linear probing.

Note- Probing:- Probing is the method in which to find an open bucket, or an element already stored in the underlying array of a hash table.

the following formula to implement  
By using linear probing

$$h(x) = (h(x) + j) \bmod \text{size}.$$

$h$  is the hash function

$x$  is the value to be stored.

$j$  is iteration of the probe. (it varies from 0 to K)

Let us consider the following data, store into the hash table of size 13 using linear probing.

0	1	2	3	4	5	6	7	8	9	10	11	12
40	27	16	28					9	75			

hash table size is 13 then insert the following elements: 16, 40, 27, 9, 75, 28

$$\text{input } 16 \rightarrow h(16) = 16 \% 13 - 3$$

$$\text{input } 40 \rightarrow h(40) = 40 \% 13 - 1$$

$$\text{input } 27 \rightarrow h(27) = 27 \% 13 = 1$$

collision occur

so we need to probe for empty location  
linear probing is

$$\begin{aligned} j &= 1 \\ h(27) &= f(27) + j \\ &= 1 + 1 = 2 \end{aligned}$$

insert 27 at 2<sup>nd</sup> position.

$$\text{input } 9 \rightarrow h(9) = 9 \% 13 = 9$$

$$\text{input } 75 \rightarrow h(75) = 75 \% 13 = 10$$

$$\text{insert}(28) = 28 \cdot 1 \cdot 13 = 2$$

Collision occurred  $\rightarrow j = 1$

$$\text{insert}(28)_1 = (28 \cdot 1 \cdot 13) + 1$$

$$= 2 + 1$$

$$= \underline{\underline{3}}$$

Third location also ~~is stated~~ contains other value so find next empty location.

$$\begin{aligned} j &= 2 \\ \text{insert}(28)_2 &= f(28) + 2 \\ &= f(28) + j \\ &= 2 + 2 \\ &= \underline{\underline{4}} \end{aligned}$$

Fourth location is empty  $ht[u] = 28$

Search for 28:-

find identify  $f(28)$  is 2

$$ht[2] \neq 28$$

find the probe

$$f(28) + 1 = 2 + 1 = 3$$

$ht[3] \neq 28$  go for next iteration.

$$f(28) + 2 = 2 + 2 = 4$$

$ht[4] = 28$  Element not found.

Note: if the location is empty simply you can say element not found.

Quadratic probing:- Quadratic probing is an open addressing method to resolve the collision in hash table. This method is eliminating the primary cluster problem which is created in linear probing.

This technique works by considering of original hash index and adding successive value of an arbitrary quadratic polynomial until the empty location is found.

In linear probing we follow

$$f(x) + 0, f(x) + 1, \dots, f(x) + k.$$

In quadratic probing we follow

$$f(x) + 0^2, f(x) + 1^2, f(x) + 2^2, \dots, f(x) + k^2$$

So

we can write the equation

$$f_i(x) = (f_i(x) + i^2) \% \text{ Sh}$$

Example:-

Let us consider the hash table size 13  
then insert the following data 16, 40, 27, 9, 75, 28

	40	27	16		28		9	75		
0	1	2	3	4	5	6	7	8	9	10

$$\text{insert}(16) = 16 \cdot 1 \cdot 13 = 3$$

$$\text{insert}(40) = 40 \cdot 1 \cdot 13 = 1$$

$$\text{insert}(27) = 27 \cdot 1 \cdot 13 = 1$$

Collision occurred. Then apply quadratic probing.

$$\begin{aligned} f(x) &= (f(x) + 1^2) \% 13 \\ &= \cancel{(27 \cdot 1 \cdot 13 + 1^2)} \% 13 \\ &= (1 + 1^2) \% 13 = 2 \end{aligned}$$

$$\begin{aligned} \text{insert}(9) &= ((f(9) + 0^2)) \% 13 \\ &= (9 + 0) \% 13 = 9 \end{aligned}$$

$$\begin{aligned} \text{insert}(75) &= (f(75) + 0^2) \% 13 \\ &= (10 + 0) \% 13 \\ &= 10 \% 13 = 10 \end{aligned}$$

$$\begin{aligned} \text{insert}(28) &= (f(28) + 0^2) \% 13 \\ &= 2 \% 13 = 2 \end{aligned}$$

Collision occurred.

$$\begin{aligned} &(f(28) + 1^2) \% 13 \\ &= (2 + 1) \% 13 \\ &= 3 \% 13 = 3 \end{aligned}$$

still collision

$$\begin{aligned} &(f(28) + 2^2) \% 13 \\ &= (2 + 4) \% 13 \\ &= 6 \% 13 = 6 \end{aligned}$$

Dynamic

Double Hashing :- we make use of two hash functions. The first hash function is  $f_1(k)$  and this function takes in our key and gives out a location on the hash table. If the location is empty easily insert the value.

In case of collision we need to use secondary hash function  $f_2(k)$ , in combination with  $f_1(k)$ .

The combined function is

$$f(k, i) = (f_1(k) + i * f_2(k)) \% \text{size}$$

$i$  is non negative integer

~~K is size~~

Example let us consider the table of size 5 then insert 10, 15, 23

10	23	15	
0	1	2	3

Consider  $f_1(k)$  is  $K \times 5$

$f_2(k)$  is  $K \times 7$

$$\text{insert}(10) = 10 \times 5 = 0$$

$$\text{insert}(15) = 10 \times 5 = 0$$

Collision occurred.  $\rightarrow i=1$

$$= ((10 \times 5) + 1 + (10 \times 7)) \% 5$$

$$= (0 + 3) \% 5$$

$$= 3$$

$$\text{Input (23)} \cdot 23 \cdot 1.5 = 3$$

Collision occurred  $i=1$

$$= (23 \cdot 1.5 + 1 * 23 \times 7) \cdot 1.5$$

$$= (3 + 2) \cdot 1.5$$

$$= 5 \cdot 1.5$$

$$= \frac{0}{2}$$

again occurred the collision  $i=2$

$$= (23 \times 5 + 2 * 23 \times 7) \cdot 1.5$$

$$= (3 + 4) \cdot 1.5$$

$$= 7 \cdot 1.5 = 2$$

Rehashing:- Rehashing is the process of re-calculating the hash codes of already stored entries to move them to another bigger size hash table. When the threshold is reached/crossed.

Rehashing will be implemented hash table reaches the load factor of 0.75 or 75% of the hash table size.

Let us consider size is 17

$$\text{so load factor is } 0.75 \times 17 = 12.75 = 13.$$

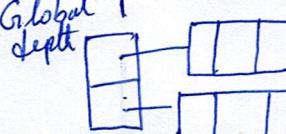
If the hash table occupancy is reached 13 then apply rehashing.

now the hash table size is ~~9+13~~ extended to next Prime number or next prime number ( $2 + \text{size}$ );

Extensible hashing:- Extensible hashing is a dynamically updatable disk-based index structure which implements a hashing technique utilizing a directory. The index is used to support exact match query. i.e. find the record with a given key compared with the B+ tree.

In this process maintain a directory of bucket address instead of just ~~hashing~~ hashing to buckets ~~search~~ directory.

The directory can grow its size power of 2. At any time, the directory consists of  $d$  levels. Let us consider directory size is 1 and Bucket size is 3

Now the hash table is  
Global depth - local depth.  


Global depth represents no. of MSB Most significant bits considered for hashing.

Let us consider the example to insert the values into hash table using extendable hashing.

Consider the initial global depth is 1 and bucket size is 3

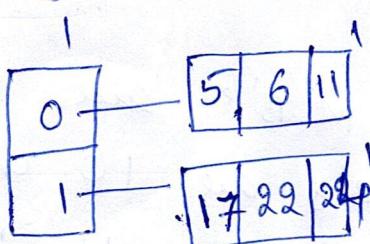
Values are 17, 5, 6, 22, 24, 11, 30, 7, 10, 21, 27.

first convert the value as binary values.

17	10001	30	11110
5	00101	7	00111
6	00110	10	01010
22	10111	21	10101
11	01011	27	11011
24	11000		

17, 5, 6, 22, 24 & 30 7 10 21 27

initially hash table is



first insert 17  
17 binary value is 10001  
its MSB is 1 so insert  
at 1's pointer bucket.

insert 5 its bucket value is 0 bcz of its MSB.

insert 6 its MSB also zero them insert at 0'

insert 22 its MSB is 1 insert after 17

insert 24 its MSB is 1 insert after 22.

insert 11 its MSB is 0 insert after 6

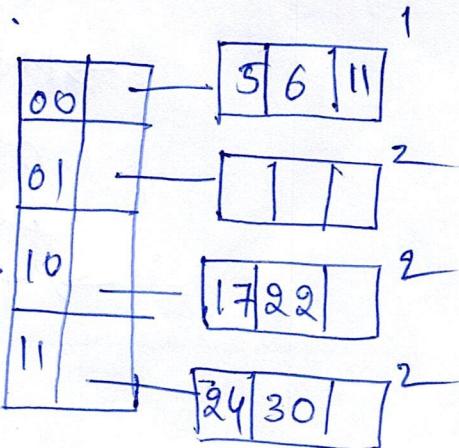
insert 30 its MSB is 1, But the

corresponding bucket is full so increase  
the directory size to power of ~~next~~ next

integer. (Bucket size is 1 it can be inserted to 2)

so power of 2 is 4.

Now the directory structure is  
like the next figure. Now you  
can split the overflow buckets.  
then apply rebalancing with  
2 MSB Bits them the table.



now insert 7 with MSB in 00 bucket as  
already full with single MSB hash value. Then  
apply rebalancing with 2 MSB's.

Now the table is

00	5	6	7
01	11	10	91
10	17	22	
11	24	30	27

Now insert 10 with MSB's are 01 insert after 00.  
Insert 21 with MSB's are 10 insert after 10.  
Insert 27 with MSB's are 11 insert after 30.

Construct the hash table with the following data.

28, 41, 19, 11, 22, 16, 12, 01, 5, 7.

28 - 11100	4 - 00100	19 - 10011	1 - 00001
22 - 10110	16 - 10000	12 - 01100	0 - 00000
05 - 00101	7 - 00111		

Advantages :- → Data retrieval is less expensive  
→ dynamic changes in hashing function.

Limitations :- → size of every bucket is fixed.  
→ Memory is wasted in pointers when the  
global depth and local depth difference  
becomes drastic.  
→ This method is complicated.