

## UNIT - I

Dictionaries : sets, dictionaries, hash tables, open hashing, closed hashing (rehashing methods), Hashing Functions (division, multiplication, universal hashing), skip lists, Analysis of skip lists.

### sets :-

A set is a collection of well defined elements.

well defined means it must be absolutely clear that which object belongs to the set & which does not.

Ex: The collection of positive numbers less than 10 is a set because given any numbers we can always find out whether that number belongs to the collection or not. But the collection of good students in your class is not a set as in the case no definite rule is supplied by the help of which you can determine whether a particular student of your class is good or not.

2. A group of singers with ages b/w 18 years & 25 years is a set because the range of ages of the singer is given so it can easily be decided that which singer is to be included & which is to be excluded.

3. A collection of red flowers is a set.

4. Collection of past presidents of the united states union is a set

5. A group of young dancers is not a set as the range of the ages of young dancers is not given & so it can't be decided that which dancer is to be considered young

### Properties of set :-

- Set is defined by capital letters.
- All the elements in the sets are enclosed within { } (curly braces)
- Every Element is separated by comma

$$\text{eg: } A = \{a, b, c, d\}$$

here  $a \in A$

$c \notin A$

### Representation of sets :-

#### 1) Statement form (Descriptive form / Describe method)

In this form the well defined descriptions of a member of a set are written and enclosed in curly braces.

for example the set of even numbers less than 15

In statement form it can be written as  
 $\{ \text{even numbers less than } 15 \}$

Roster Form :- (Tabular form / listing method)

In this form all the elements of a set are listed.

For example set of natural numbers less than 5

It can be written as  $N = \{1, 2, 3, 4\}$

Set Builder Form :-

In this representation of a set, the element of the set is described by using a symbol  $x$  or any other variable followed by colon ( $:$ ) or ' $|$ ' is used to denote such that men associate the property possessed by the elements of the set & enclose the whole description in braces.

Ex:  $A = \{2, 4, 6, 8\} \text{ s.t. } 1 \leq n \leq 4$

$A = \{x : x = 2n, n \in N \text{ & } 1 \leq n \leq 4\}$   
read as  $A$  is the set of elements  $x$  such that  $x = 2n$  &  $1 \leq n \leq 4$

Types of sets :-

Empty set: A set which does not contain any element. It is denoted by  $\{\}$  or  $\emptyset$ .

Singleton set :-

A set which contains a single element

Ex: There is only one apple in basket

## Finite set :-

A set which consists of a definite no. of elements

Ex: A set of natural number upto 5

$$A = \{1, 2, 3, 4, 5\}$$

## Infinite set :-

A set which is not finite

Ex: A set of all natural numbers.

$$A = \{1, 2, 3, 4, \dots\}$$

## Equal sets -

Two sets  $A$  &  $B$  can be equal only if each element of set  $A$  is also the element of the set  $B$ . Also if two sets are subsets of each other they are said to be equal

$$A = B$$

$$ACB \& BCA \Rightarrow A = B$$

## Equivalent sets :-

two sets are equivalent if the no. of elements in both the sets is equal and it is not necessary that they have same elements or they are subset of each other

$$P: \{1, 3, 9, 5, -7\} \& Q: \{5, -7, 3, 1, 9\}$$

then  $P$  &  $Q$  are Equal sets

$$P: \{1, -7, 10\} \& Q: \{A, B, C\}$$

$P$  &  $Q$  are Equivalent sets

### Disjoint sets :-

Two sets  $A \cup B$  are said to be disjoint if the set does not contain any common elements.

$$\text{Ex: } A = \{1, 2, 3, 4\} \quad B = \{5, 6, 7\}$$

### Subsets :-

A set  $A$  is said to be subset of  $B$  if every element of  $A$  is also an element of  $B$ , denoted as  $A \subseteq B$ , even null set is considered to be subset of any set.

$$\text{Ex: } A = \{1, 2, 3\}$$

$$\text{then } \{1, 2\} \subseteq A$$

similarly other subsets of  $A$  are  $\{\}, \{1\}, \{2\}, \{3\}$ ,  $\{1, 2\}, \{1, 3\}, \{2, 3\}, \{1, 2, 3\}$

If  $A$  is not a subset of  $B$  then it is denoted as  $A \not\subseteq B$

### Proper Subset :-

If  $A \subseteq B$  &  $A \neq B$  then  $A$  is called the proper subset of  $B$  & can be written as  $A \subset B$

Ex: If  $A = \{2, 5, 7\}$  is a subset of  $B = \{2, 5, 7\}$  then it is not a proper subset of  $B$ .

But  $A = \{2, 5\}$  is a subset of  $B = \{2, 5, 7\}$  and is a proper subset also.

## Superset:-

Set A is said to be superset of B if all the elements of set B are elements of set A it is represented as  $A \supset B$

Ex.  $A = \{1, 2, 3, 4\}$  & set  $B = \{1, 3, 4\}$   
then A is superset of B.

## operations on sets:-

### union of sets:-

If set A & B are two sets then  $A \cup B$  is set that contains all elements of set A and set B

$$A = \{1, 2, 3\} \quad B = \{4, 5, 6\}$$

$$A \cup B = \{1, 2, 3, 4, 5, 6\}$$



### Intersection of sets:-

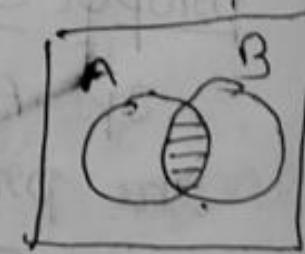
If A & B are two sets then  $A \cap B$  is the set that contains only common elements b/w A & set B

$$\text{1) } A = \{1, 2, 3\} \quad B = \{3, 4, 5\}$$

$$A \cap B = \{3\}$$

$$\text{2) } A = \{1, 2, 3\} \quad B = \{4, 5, 6\}$$

$$A \cap B = \emptyset \text{ or } \{\}$$



### complement of sets:-

complement of any set P, is the set of all elements in universal set that are not in set P. It is denoted by  $P'$

$$P' = U - P$$



the final key-value pair. So if we insert duplicate keys then our original data is lost.

3. The key attribute is case sensitive hence key is not the same as KEY

4. The key attribute is immutable hence we cannot have an array or list as keys as we can only have numbers, strings, etc. as the key.

### Operations in dictionary :-

Add or Insert :- In this operation a new pair of keys and values is added in the dictionary or associative array object.

### Replace or reassign :-

In the replace or reassign operation the already existing value that is associated with a key is changed or modified. In other words a new value is mapped to an already existing key.

### Delete or remove :-

In the delete or remove operation the already present elements are unmapped from the dictionary or associative array object.

the final key-value pair. So if we insert duplicate keys then our original data is lost.

3. The key attribute is case sensitive hence key is not the same as KEY

4. The key attribute is immutable hence we cannot have an array or list as key.

• Keys can only have numbers, strings, etc. as the key.

### Operations in dictionary :-

#### Add or Insert :-

This operation adds a new pair of keys and values if added in the dictionary or associative array object.

#### Replace or reassign :-

In the replace or reassign operation the already existing value that is associated with a key is changed or modified. In other words a new value is mapped to an already existing key.

#### Delete or remove :-

In the delete or remove operation the already present element is unmapped from the dictionary or associative array object.

Find on lookup: In the find operation the value associated with a key is searched by passing the key as a search argument.

insertItem(5, A) {5, A}  
insertItem(7, B) {5, A}(7, B)  
insertItem(2, C) {5, A}, {7, B}, {2, C}  
insertItem(8, D) {5, A}, {7, B}, {2, C}, {8, D}  
insertItem(2, E) {5, A}, {7, B}, {2, E}, {8, D}

findItem(7)

findItem(4)

removeItem(5) {7, B}(2, E), {8, D}

replaceItem(7, M) {7, M}(2, E), {8, D}

### Hash tables:

A hash table is one of the most important data structures that uses a special function known as hash function that maps a given value with a key to access the elements faster.

A hash table is DS that stores some information & the information has basically two main components i.e key & value, the hash table efficiency of mapping depends upon the efficiency of the hash function used for mapping.

## Drawback of hash function

A hash function assigns each value with a unique key. Sometimes hash table uses an imperfect hash function that causes a collision because the hash function generates the same key of two different values.

## Hashing :-

Hashing is one of the searching techniques that uses a constant time complexity in hashing is  $O(1)$ . We know linear search worst time complexity is  $O(n)$  & binary search worst case time complexity is  $O(\log n)$ . In both searching techniques the searching depends upon the no. of elements but we want the technique that takes a constant time. So hashing technique came that provides a constant time.

The main idea behind hashing is to create the (key / value) pairs. If the key is given, then the algorithm computes the index at which the value would be stored. It can be written as:

Index, e.g.  $\text{hash}(\text{key})$

There are three ways of calculating hash function

1) Division method

2) Folding method

3) Mid Square method

4) Multiplication method

## hashing functions

### Division method

In this method the hash function

can be defined as

$$h(K) = K \% m$$

where  $m$  is size of hash table

generally  $m$  is chosen to prime number.

#### Example

$$K = 88$$

$$M = 11$$

$$h(88) = 88 \% 11$$

$$K = 78$$

$$M = 11$$

$$h(78) = 78 \% 11$$

$$K = 50$$

$$M = 11$$

$$h(50) = 50 \% 11$$

$$= 6$$

6	88
1	78
2	
3	
4	
5	
6	50
7	
8	
9	
10	

Pros:-

- This method is quite good for any value of  $M$ .
- The division method is very fast since it requires only a single division operation.

Cons:-

Sometimes extra care should be taken to choose the value of  $M$ .

### Mid Square Method

It is very good hashing method. It involves two steps to compute hash value.

1. Square the value of the key  $K$ , i.e.  $K^2$ .
2. Extract the middle  $\sigma$  digits as the hash value.

### Formula:

$$h(K) = \text{sh}(K \times K)$$

The value of  $\sigma$  can be decided based on the size of the table.

Ex:-

Suppose hash table has 100 memory locations so  $\sigma = 2$  because two digits are required to map the key to the memory location.

$$K = 60$$

$$\begin{aligned} K \times K &= 60 \times 60 \\ &= 3600 \end{aligned}$$

$$h(60) = 60$$

## Digit Sonding method

This method involves two steps

1. Divide the key-value  $k$  in to a no. of parts i.e.  $k_1, k_2, k_3, \dots, k_n$  where each pair has the same no. of digits except for the last part man can have lesser digits than other parts.

2. Add the individual parts. The hash value is obtained by ignoring the last carry if any.

Ex: 12345

$$k_1 = 12, k_2 = 34, k_3 = 5$$

$$S = k_1 + k_2 + k_3$$

$$= 12 + 34 + 5 = 5019$$

$$n(k) = 5$$

$$S = 5019$$

## Multiplication method

choose a constant value  $A$  such that  $0 < A < 1$

1. Multiply the key value with  $A$

2. Extract the fractional part of  $ka$

3. Extract the result of above step by

4. Multiply the result of above step by the size of the hash table i.e.  $H$

5. The resulting hash value is obtained

6. The resulting hash value is obtained by taking the floor of result obtained in step 4

Formula:  $n \equiv f_{1001}(\text{MCICA mod } 1)$

eg. 2 out of 1001 boxes will

be of size  $M = 8^{\text{size}}$

using as  $K = 1001 - 1001 \times 0.618033$

$A = \text{constant} = 3.141592$

Wipib 2000 mm will be stored in a suggestion box.

$R = 23$  total no. of 19000

$M = 1000$  max wipib

$$n(23) = f_{1001}(1000(23 \times 0.618033)) \% .10$$

$$\begin{aligned} &= f_{1001}(142.143 \% .10) \\ &= 142 \% .10 \end{aligned}$$

= 2

$$n(23) = f_{1001}(1000(23 \times 0.618033 \bmod 1))$$

$$f_{1001}(1000(14.2147 \% .1))$$

$$f_{1001}(1000(0.2147))$$

$$f_{1001}(2.147) = 2$$

UNIVERSAL HASHING → last page of UNIT-1

(open hashing :- (closed addressing)) \*

collision handling techniques

If two pieces of data share the same value in a hash table then it is known as collision in hashing. As the name suggests there collisions are handled

using different techniques known as collision handling techniques. These techniques are also known as collision

resolution techniques.

The techniques that are used to handle or resolve the collisions are basically classified into two types they are:

1) open hashing (or) separate chaining (or )  
closed chaining.

2) closed hashing (or) open Addressing

### Open Hashing:-

The first collision resolution or handling technique, "open hashing", is popularly known as separate chaining. This is a technique which is used to implement an array as a linked list known as a chain. It is one of the most used techniques by programmers to handle collisions.

Basically a linked list data structure is used to implement the separate chaining technique.

When a number of elements are hashed in to the index of singly linked list then they are inserted in to a singly linked list. This singly linked list is the linked list which we refer to as a chain in the open hashing technique.

## Advantages of open hashing:

- this method is simple to implement & understand
- it can be implemented when we do not know how frequently keys will be inserted or deleted.
- open hashing is less sensitive to the load factors or hash function.

## Disadvantages:

- The cache performance of separate chaining method is poor as the keys are stored using a singly linked list.
- A lot of storage space is wasted as some parts of the hash table are never used.
- In the worst case the search time can become  $O(n)$ , this happens only if the chain becomes too dense.
- Extra storage is used for storing pointers of the chains.

## Functions used in open hashing

- 1) get CK key :- If the key is present in hash table, get CK key gives the value corresponding to the key (hash table).

2. getSize() :- It will return the size of hash table.

3. add() :- changes an existing valid key-value pair in hash table if it already exists before adding a new one.

4. remove() :-

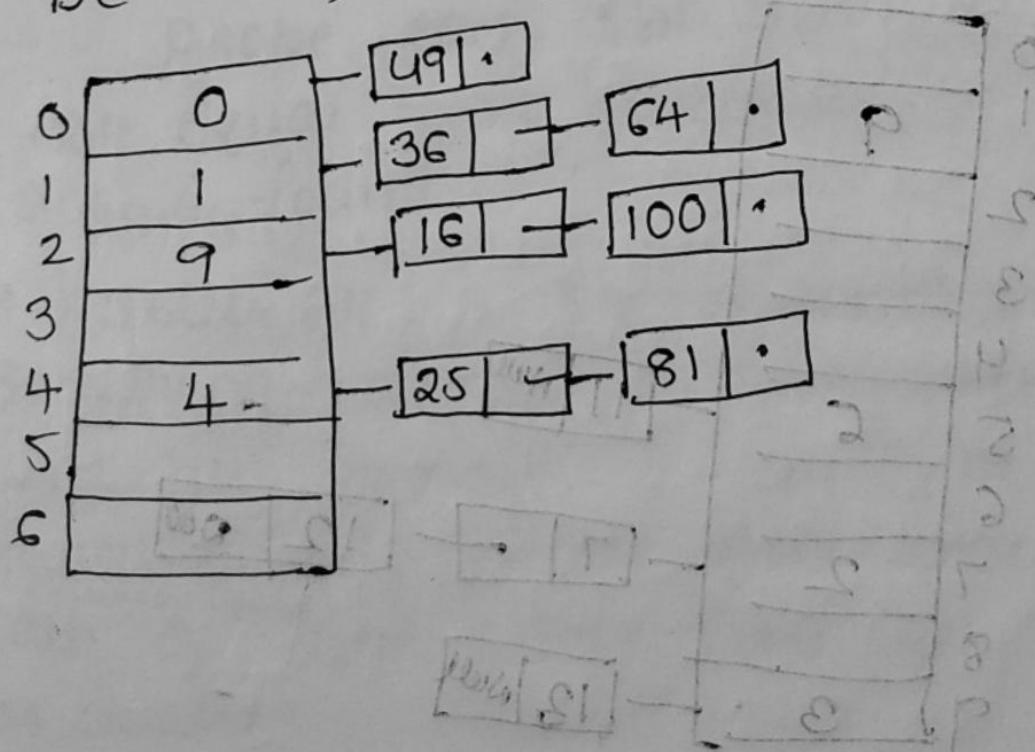
The remove function will remove both the key & value pair.

5. isEmpty() :-

It returns the value "true", if the size is zero.

Example :-

Consider the keys 0, 1, 4, 9, 16, 25, 36, 49, 64, 81, 100 let the hash function be  $h(x) = x \% 7$

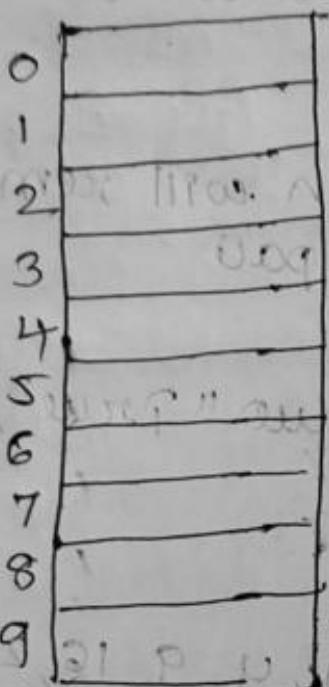


Example: 2

A: 3, 2, 9, 6, 11, 13, 7, 12

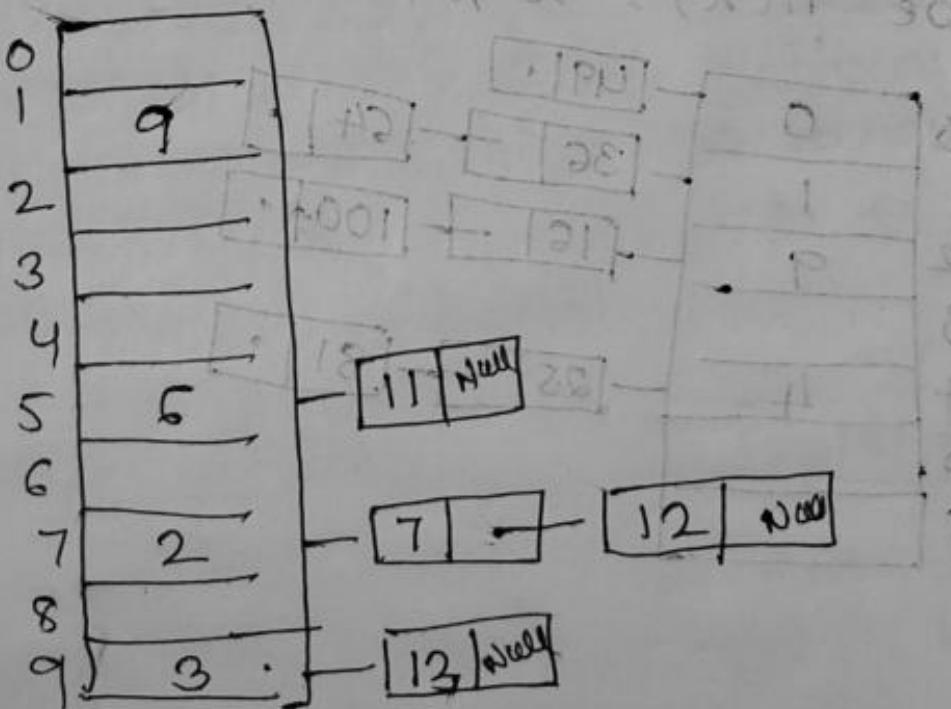
$$n(k) = 2k+3, m=10$$

use division method



key	location
3	$[2(3)+3] \% 10 = 9$
2	$[(2 \times 2) + 3] \% 10 = 7$
9	$[(2 \times 9) + 3] \% 10 = 1$
6	$[(2 \times 6) + 3] \% 10 = 5$
11	$[(2 \times 11) + 3] \% 10 = 5$
13	$[(2 \times 13) + 3] \% 10 = 9$
7	$[(2 \times 7) + 3] \% 10 = 7$
12	$[(2 \times 12) + 3] \% 10 = 7$

result:



## Closed Hashing or Open Addressing:-

The second most collision resolution technique, closed hashing is a way of dealing with collisions, similar to separate chaining process. In open addressing, the hash table alone stores all of its elements. The size of the table should always be greater than or equal to the total no. of keys at all times.

functions used in closed hashing:-

1. Insert ( $K$ ) :- Up till a space is left unfilled keep probing. place the key " $K$ " in the first empty slot you find.
2. Search ( $K$ ) :- probe each slot until the key is not equal to  $K$  or until an empty slot is found.
3. Delete ( $K$ ) :- It's interesting to delete something. The search may fail when we just remove a key & then perform search operation. the slots that are a part of deleted key slots are considered as deleted.

Several techniques to perform implementation of closed hashing:-

1) Linear probing :- Searching for Empty slot.

In linear probing, the hash table undergoes clear examination, starting from the hash's initial or beginning point of the slot that is obtained after the calculation is already occupied then we should look for different one.

The function that is responsible for performing rehashing is "rehash  $(n+1) \% \text{table\_size}$ "

Difficulties faced with linear probing :-

→ primary clustering :-

It is one of the major issues that are caused with the linear probing technique. Many elements that are consecutive to each other generally form clusters or a group of scattering which in turn makes the hash table more difficult to find an empty slot or search for an element.

## → Secondary clustering :-

Secondary clustering is not as severe as primary clustering & the elements or records that must be placed within the same location are only allowed to share a collision chain which is also known as probe sequence, if they begin at the same location.

→ clustering is the only problem in linear probing. If clustering can be reduced within this mechanism, then this can be considered one of the best collision resolution technique.

## 2. Quadratic probing :-

The function used for rehashing is as follows:

$$\text{rehash(key)} = (n+1) \% \text{table-size}$$

Ex: let us consider a simple hash function as "key mod 7" & sequence of keys as 50, 700, 76, 85, 92, 73, 101

Sol:

key 50

$$h(x) = x \% 7$$

$$50 \% 7$$

$$= 1$$

50 is inserted at index 1

key 700

$$h(x) = x \% 7$$

$$700 \% 7$$

$$= 0$$

700 is inserted at index 0

key 76

$$h(x) = x \% 7 \text{ parallel hashing}$$

$$= 76 \% 7 \rightarrow \text{parallel hashing}$$

$$= 6 \rightarrow \text{index 6}$$

76 is inserted at index 6

76 is inserted at index 6

key 85

$$h(x) = x \% 7 \text{ will be parallel hashing}$$

$$= 85 \% 7$$

$$= 1$$

collision occurs.

0	700
1	50
2	85
3	92
4	73
5	101
6	76

$$\text{rehash}(x) = (h(x) + 1) \% S$$

$$\text{rehash}(x) = (1 + 1) \% 7$$

$$= 2$$

$$\text{key} = 92 \rightarrow h(92) = (\text{first digit})$$

$$h(x) = x \% 7$$

$$= 92 \% 7$$

$$= 1$$

101 collision occurs

$$\text{rehash}(x) = (h(x) + 1) \% S$$

$$= (1 + 1) \% 7 = 2$$

collision

$$\text{rehash}(x) = (h(x) + 2) \% S$$

$$= (1 + 2) \% 7$$

$$= 3$$

92 is inserted at index 3

$$key = 73$$

$$h(x) = x \% 7 \quad (\text{and result is } 101 \% 7 = 4)$$

$$273 \% 7 \rightarrow 4 \quad (\text{and result is } 101 \% 7 = 4)$$

but only  $38 \% 7 = (101 + 100) \% 7 = 101 \% 7 = 4$   
 collision

$$\text{rehash}(x) = (h(x) + 1) \% S$$

$$21 \% (8+8 + 100) \% 7 = 101 \% 7 = 4$$

will insert 73 at index 4 next.

$$\text{key} = 101$$

$$h(x) = x \% 7$$

$$= 101 \% 7$$

3 (collision)

$$\text{rehash}(x) = (h(x) + 1) \% S$$

$$101 \% 7 \rightarrow 4 \quad (\text{collision})$$

$$\text{rehash}(x) = (h(x) + 2) \% S$$

$$= (3+2) \% 7$$

insert 101 at index 5

## 2. Quadratic probing

Quadratic probing is an open-addressing scheme where we do one iteration if the given hash value  $x$  collides in the hash table.

Let  $\text{hash}(x)$  be the slot index computed using hash function

- If the slot  $\text{hash}(x) \% s$  is full, then we try  $(\text{hash}(x) + 1 * 1) \% s$
- If  $(\text{hash}(x) + 1 * 1) \% s$  is also full, then we try  $(\text{hash}(x) + 2 * 2) \% s$
- If  $(\text{hash}(x) + 2 * 2) \% s$  is also full then we try  $(\text{hash}(x) + 3 * 3) \% s$
- This process is repeated for all  $m$  values of  $i$  until an empty slot is found.

Ex :-

Let us consider a simple hash function as "key mod 7" and sequence of keys as 50, 700, 76, 85, 92, 73, 101

$$\rightarrow K = 50$$

$$n(x) > K \% 7$$

$$= 50 \% 7$$

$$= 1$$

$$\rightarrow K = 700$$

$$n(x) > 700 \% 7$$

$$\text{Index} = 0$$

$$\rightarrow K = 76$$

$$n(x) = 76 \% 7$$

$$= 6$$

$$\rightarrow K = 85$$

$$n(x) = 85 \% 7$$

$\rightarrow 1$  (collision)

0	: 700
1	: 50
2	: 85
3	: 73
4	: 101
5	: 92
6	: 76

$$(\text{hash}(x) + 1 * 1) \% 7$$

$$(1+1) \% 7$$

so as to reduce the number of collisions  
in bucket B = 2

$$\rightarrow K = 92 \quad n(x) = 92 \% 7$$

= 1 (collision)

$$(\text{hash}(x) + 1 * 1) \% 7 \text{ no collision}$$

$$(1+1) \% 7 = 2 \text{ (collisions)}$$

$$(\text{hash}(x) + 2 * 2) \% 7 \text{ no collision}$$

$$(1+4) \% 7 = 5 \text{ (collisions)}$$

92 is inserted at index 5

$$\rightarrow K = 73$$

$$n(x) = 73 \% 7 = 3 \text{ (at index 3)}$$

73 is inserted at index 3

$$\rightarrow K = 101$$

$$n(x) = 101 \% 7 = 4 \text{ (collisions)}$$

101 is inserted at index 4

$$(\text{hash}(x) + 1 * 1) \% 7$$

$$(3+1) \% 7 = 4 \text{ (max)}$$

$$4 \% 7 = 4$$

101 is inserted at index 4

101 is inserted at index 4

101 is inserted at index 4

### 3. Double hashing:

Double hashing is worse on a similar idea to linear & quadratic probing. use a big table & hash into it. whenever a collision occurs, choose another spot in table to put the value. The difference here is that instead of choosing next opening, a second hash function is used to determine the location of the next spot. For example given hash function  $H_1$  &  $H_2$  & key.

- check location  $hash_1(key)$ . If it is empty put record in it
- If it is not empty calculate  $\& hash_2(key)$
- check if  $(hash_1(key) + hash_2(key)) \mod M$  is open, if it is, put it in
- repeat with  $hash_1(key) + 2hash_2(key)$ ,  $hash_1(key) + 3hash_2(key)$  & so on until an opening is found.

Example :-

Insert the keys 27, 43, 69, 72 into hash table of size 7 where first hash function is  $h_1(k) = k \mod 7$  & second hash-function is  $h_2(k) = 1 + (k \mod 5)$ .

Step 1 : Insert 27

$27 \% 7 = 6$ , location 6

a) Empty so insert 27

Value to 6 8101

Step 2 : Insert 43

$43 \% 7 = 1$ , location 1

a) Empty so insert 43 in to 1 8101

Step 3 : Insert 692

$692 \% 7 = 6$ , location 6 is already

being occupied this is a collision

so we need to resolve this collision  
using double hashing

$$h_{\text{new}} = [h_1(692) + i * (h_2(692)) \% 7]$$

$$[6 + 1 * (1 + 692 \% 5)] \% 7$$

$$= \frac{9 \% 7}{2}$$

Now 2 is empty 8101

so we can insert 692 into 2<sup>nd</sup> 8101

Step 4 : Insert 72

$72 \% 7 = 2$ , but location 2 is  
already being occupied this is a  
collision

$$h_{\text{new}} = [h_1(72) + i * (h_2(72)) \% 7]$$

$$[2 + 1 * (1 + 72 \% 5)] \% 7$$

$$= \frac{5 \% 7}{5}$$

0	
1	43
2	692
3	
4	
5	72
6	27

we can insert 72 in to 5m slot

2nd Method :-

The following function denotes double hashing

$$(H_1(K) + i * H_2(K)) \% \text{table size}$$

where  $H_1(K)$  &  $H_2(K)$  are 2 hash

functions & table size is size of hash table.

A typical first hash function is given by:

$$H_1(K) = K \% \text{table size}.$$

A typical second hash function is given by:

$$H_2(K) = CP - (K \% CP)$$

where  $CP$  is a no. that must be smaller than the size of the hash-table & should be coprime with size of table.

Ex:-

let us consider keys 20, 34, 45, 70, 56

let size of table be 11 so the first hash function will be given by:

$$H_1(K) = K \% 11$$

for the second hash function,  $CP$  should be chosen in such a way  $\text{NCF}(CP, 11) = 1$

second hash function cp showed be  
chosen  $H2(K) = 8 - (K \% 8)$

key = 20

$h1(20) = 20 \% 11 = 9$  (No collision)  
20 is inserted at 9

$\rightarrow \text{key} = 34$

$h1(34) = 34 \% 11 = 1$

34 is inserted at index 1

$\rightarrow \text{key} = 45$

$h1(45) = 45 \% 11 = 1$   
(Collision)

$h2(45) = 8 - (45 \% 8) = 3$

$(1+1+3)\%11 = 4$

45 is inserted at index 4

$\rightarrow \text{key} = 70$

$h1(70) = 70 \% 11 = 4$

$h2(70) = 8 - (70 \% 8) = 2$

$(4+1+2)\%11 = 6$

70 is inserted at index 6

$\rightarrow \text{key} = 56$

$h1(56) = 56 \% 11 = 1$  (Collision)

$h2(56) = 8 - (56 \% 8) = 8$

0	
1	34
2	
3	56
4	45
5	
6	70
7	
8	
9	20
10	

$$(1+1*8)\%11 = 9 \text{ (collision)}$$

$$(1+8*8)\%11 = 6 \text{ (collision)}$$

$$(1+3*8)\%11 = 3$$

so it is inserted at index 3.

## SKIP LIST:

Skip List 'S' consists of a series of lists  $\{S_0, S_1, S_2, \dots, S_h\}$  such that each list  $S_i$  stores a subset of the items sorted by increasing key.

→  $h$  represents the height of the skip list.

5

→ Each list has two sentinel nodes  $+\infty$  &  $-\infty$ .  $-\infty$  is smaller than any possible key in the list &  $+\infty$  is greater than any possible key in the list. Hence the node with  $-\infty$  key is on left most position & node with  $+\infty$  key is always the right most node in the list.

→ For visualization, it is customary to have the  $S_0$  list on the bottom & lists  $S_1, S_2, \dots, S_h$  above it.

→ Each node in a list has to be "sitting" on another node with same key below it meaning that if  $L_{i+1}$  has a node

with key K then Li, Lj, ... all valid  
i.e. below it will have same key in mem  
→ skip list also needs to maintain the  
head pointer i.e. reference to the  
first member (sentinel node - ob) in the  
top most list  
→ on average the skip list will have  
 $O(n)$  space complexity

### Ideal Skip list :-

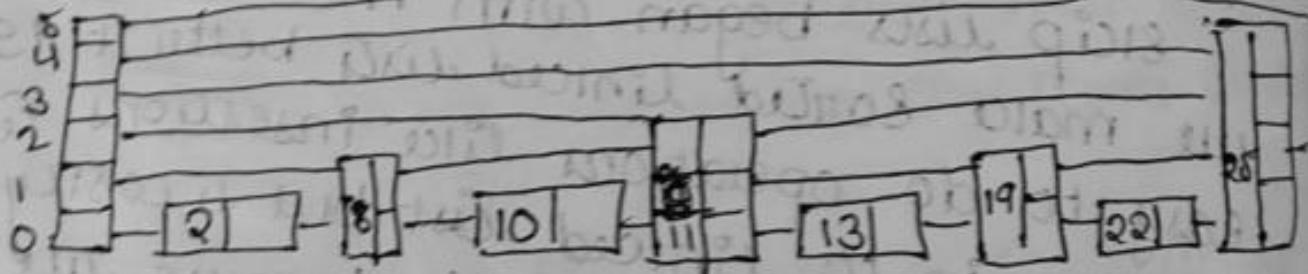
skip lists began with the idea how can  
we make sorted linked lists better if it is  
easy to do operations like insertions &  
deletion in to linked list but is costly  
to locate items efficiently because we  
have to walk through the list one item at  
a time. If we could over multiple items  
at a time however then we could  
perform searches efficiently.

Intuitively, a skip list is a data structure  
that encodes a collection of sorted linked  
lists, where linked skip over 2, then 4, then  
8, & so on, elements with each link.

To make this more concrete, imagine  
a linked list, sorted by key value. There  
are two nodes at either end of the list -  
called head & tail. trace every other  
entry of this linked list & extend it  
up to a new linked list & extend it.

with  $\frac{1}{2}$  as many entries. Now take every other entry of this linked list & extend it up to another linked with  $\frac{1}{4}$  as many entries as original list & continue this until no elements remain. The head & tail nodes are always kept clearly we can repeat this process  $\lceil \log n \rceil$  times.

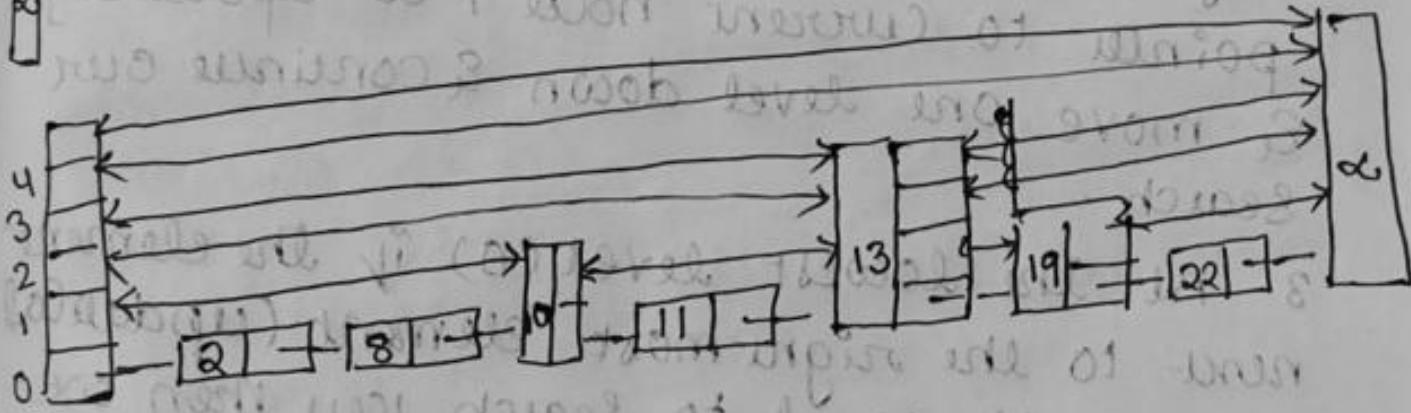
The result is something that we will call the ideal skip list.



### Randomized Skip list :-

Unfortunately since a perfectly balanced binary tree, the ideal skip list is too pure to able to use for a dynamic data structure. As soon as a single node was added to the middle of its lists all the heights following it would need to be modified, but we can relax this requirement to achieve an efficient data structure. In the ideal skip list every other node from level  $i$  is extended up to level  $i+1$ . Instead how about we do it randomly?

Suppose that we have built a skip list up to some level  $i$ , & we want to extend this to level  $i+1$ . Imagine a node at level  $i$  tossing a coin. If the coin comes up heads (with probability  $1/2$ ) the node promotes itself to level  $i+1$ , otherwise it stops here. By the nature of randomization the expected no. of nodes at level  $i+1$  will be half the no. of nodes at level  $i$ . Thus the expected no. of nodes at level  $K$  will be  $n/2^K$ .



### operations on skip list :-

#### 1. Insertion operation :-

To insert any element in a list.

#### 2. Search operation :-

To search any element in a list.

#### 3. Deletion operation :-

To delete any element from a list.

## Searching operation

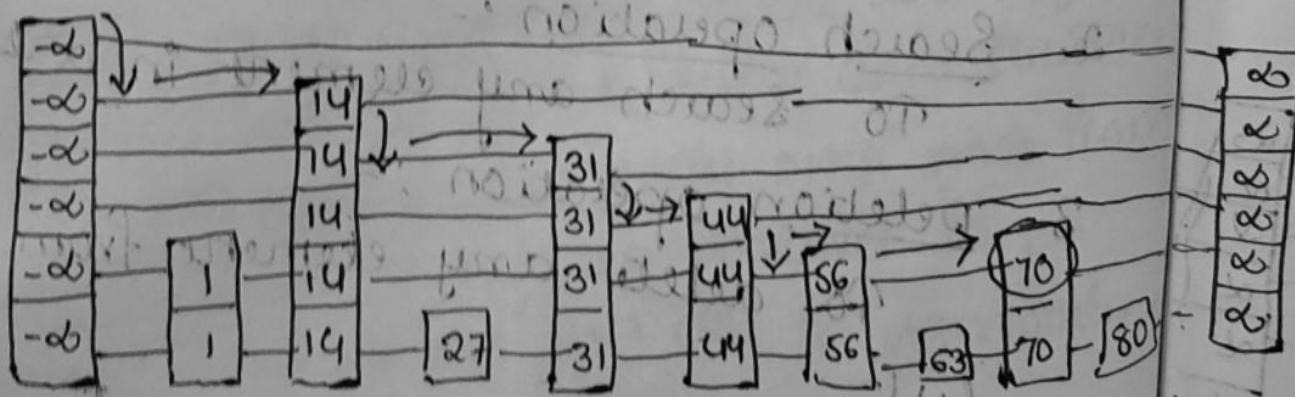
To searching an element, is very simple to approach for searching a spot for inserting an element in a skip list.

### Algorithm

1. If key or next node is less than search key then we keep on moving forward on the same level.
2. If key or next node is greater than the key to be inserted then we ignore the pointer to current node i at update[i] & move one level down & continue our search.
3. At the lowest level (0) if the element next to the right-most element (update[0]) has key equal to search key then we have found key otherwise failure.

Example :-

for finding 70



we will simply start with the left most node in the list on top. Go down if the next value in the same list is greater than the one we are looking for, go right otherwise.

Inserting in skip list :-  
To insert a new node we would first find the location where this new key should be inserted in the list at the very bottom. This can be simply done by reusing the search logic, we start traversing to the right from the list on top & we go one step down if the next item is bigger than the key that we want to insert, else go right. Once we reach at a position in bottom list where we can't go more to the right, we insert the new value on the right side of that position.

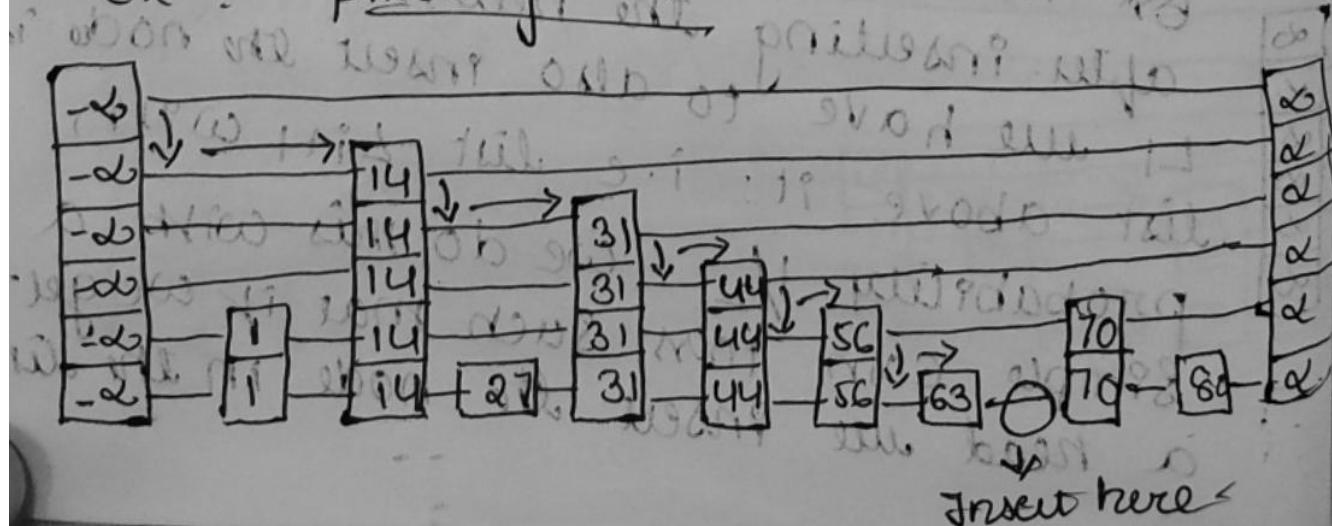
Now comes the probabilistic part after inserting the new node in the list we have to also insert the node in list above it. i.e. list  $L_{i+1}$  with probability  $1/2$ . We do this with a simple going fast, such that if we get a head we insert the node in the list  $L_i$ .

00
01
02
03
04
05

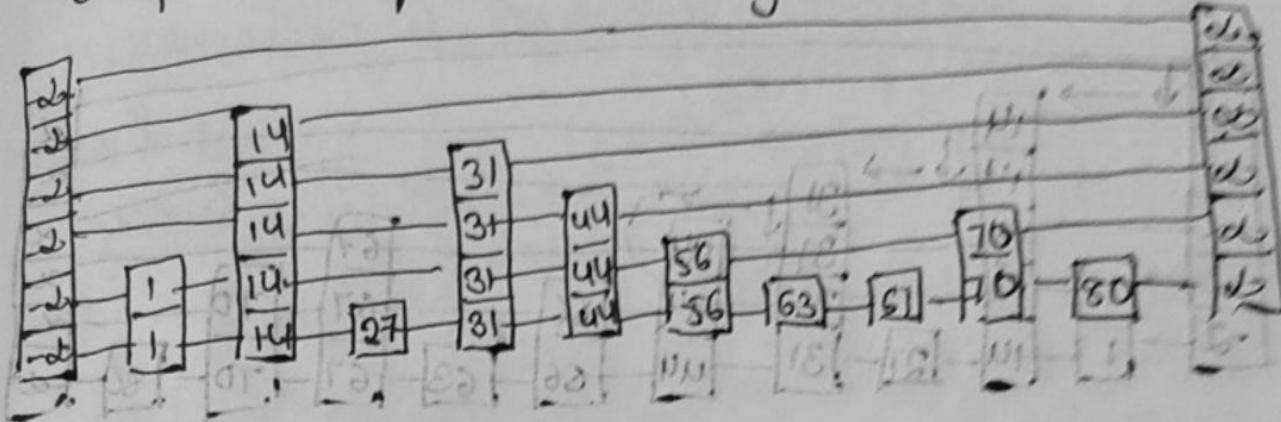
Let's toss the coin again, we stop when we get a tail else may repeat the coin toss & still we keep getting heads, even if we have to insert a new layer (list) at the top.

So we also need to re-hook the links for the newly inserted node; the new node needs to have its below reference pointing to the node below and the node below would have above reference pointer to the new node, we simply traverse toward left from the node below if & return above pointer from a node for which above pointer is not null. The right reference of the left neighbor is used to update the right reference of the new node. Also the left reference for the new nodes right neighbor is also updated to the new node. This re-linking operation is actually pretty easy to implement with linked lists.

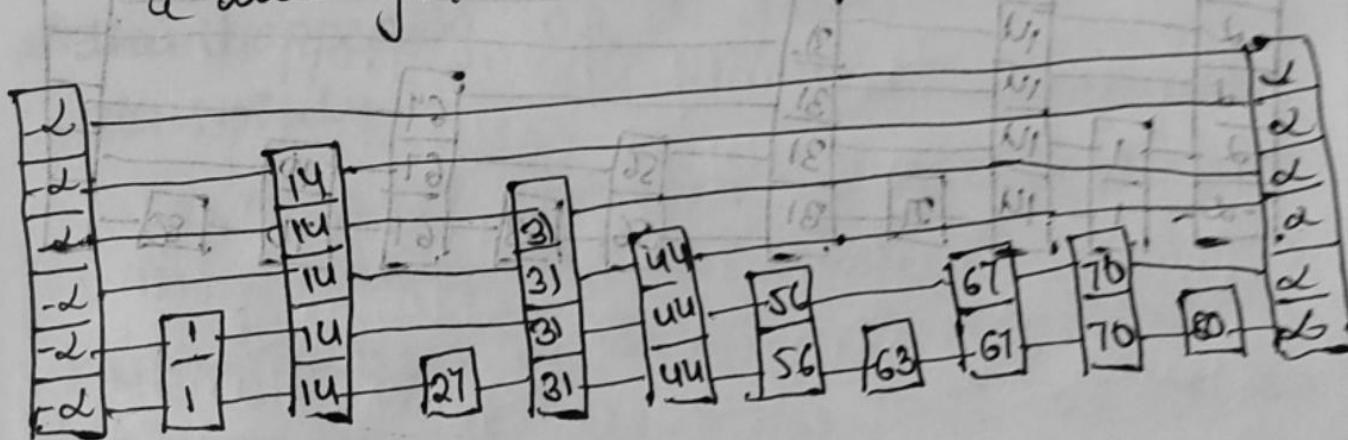
Ex: Finding, 67



skip list after inserting 67

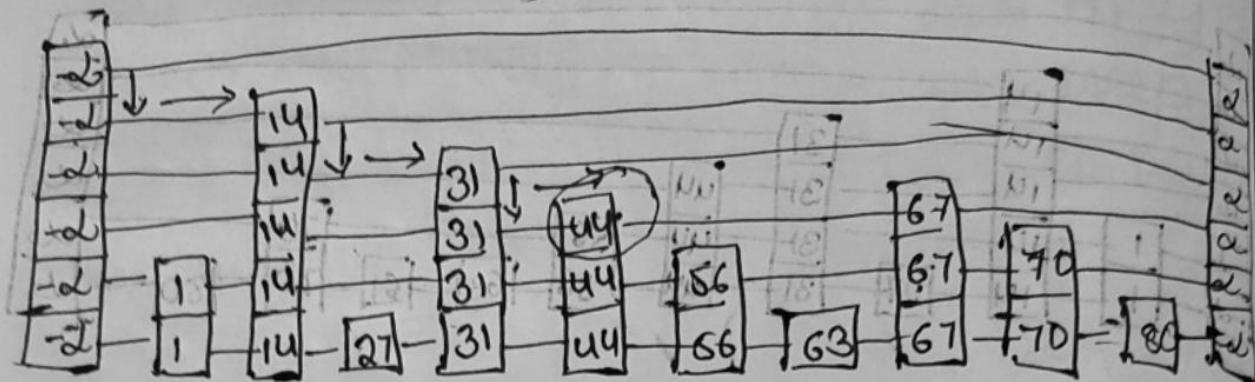


skip list after re-linking  
Suppose we got head while tossing the coin.  
Again we got head while tossing,  
& we got tail for the third time.

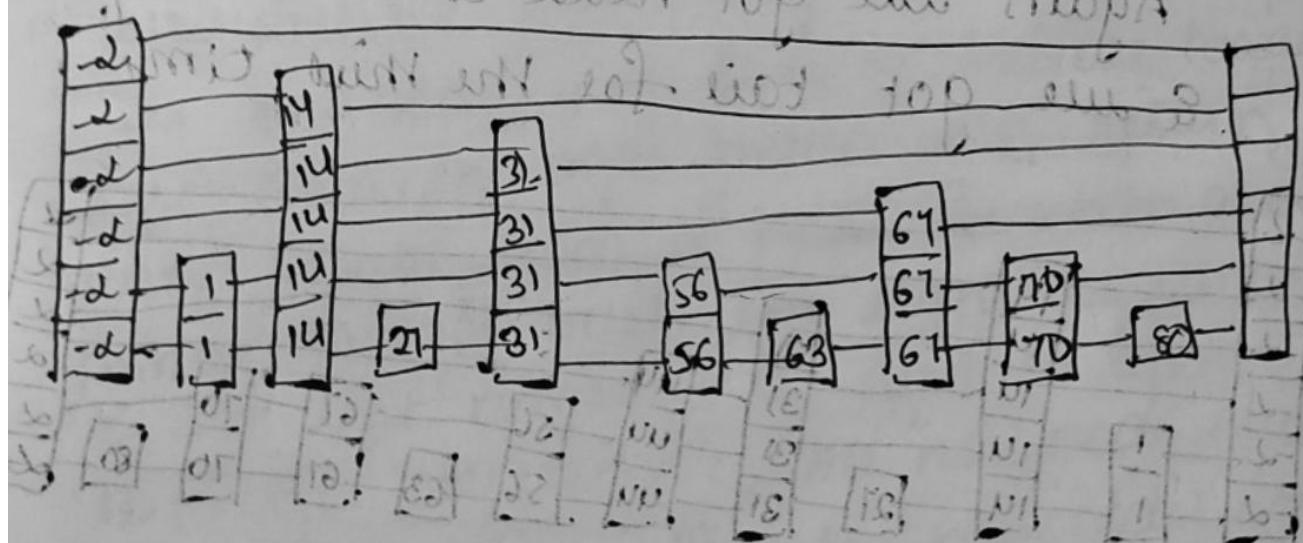


Deletion in the skip list :-  
we first perform search operation  
to find the location of the node. If we find the node, we simply delete it at all levels. Before deleting a node, we simply ensure that no other node is pointing to it.

Ex Delete 44



After finding & releasing all references to this node & move to node below & release that node as well



### Analysis of SICP List

- Search complexity is  $O(\log n)$  in average case &  $O(n)$  in worst case
- Delete complexity is  $O(\log n)$  in average case &  $O(n)$  in worst case
- Insert complexity is  $O(\log n)$  in average case &  $O(n)$  in worst case
- Space complexity is  $O(n \log n)$  in worst case

→ The height of a skip list containing elements in at most  $\log n^2$

universal hashing :-

If we have a fixed hash function such that all keys uses this single hash function, then we may be unfortunate that the keys selected may be such that they all hash to the same slot in hash table. This would be result in an average search time of  $O(n)$ . We want to avoid this happening or at least limit its possibility.

Universal hashing is a way to avoid this situation. In Universal hashing the hash function is selected randomly independent of the keys that are to be stored.

In universal hashing at the beginning of the execution we choose a hash function randomly from a carefully designed family of functions. This guarantees that no single input would result in worst case situation. Since the hash functions are random the behaviour would be different for

each execution which would yield average case performance on each input.

Let ' $H$ ' be a finite set of hash functions. Let 'U' be a given universe of keys that map a given universe 'V' of keys into range  $\{0, 1, 2, \dots, m-1\}$ , then a set of functions is said to be universal if for each pair of distinct keys  $K, L \in U$  the no. of hash function  $h \in H$  for which  $h(K) = h(L)$  is at most  $1/m$ .

This means that w.p.m a hash function randomly chosen from the set 'H' the possibility of collision between two distinct keys  $K, L$  is at most the chance of  $1/m$  i.e. a collision occurring if  $h(K) = h(L)$  were chosen randomly & independently from the set  $\{0, 1, 2, \dots, m-1\}$ .

A most plausible notion about a hash function is that it maps a large space of inputs into a much smaller space of outputs. In other words, it is a many-to-one mapping.