

3. Hashing

Concept:

Hashing is an effective way to reduce the number of comparisons. Actually hashing deals with the idea of proving the direct address of the record where the record is likely to store.

e.g.

key record hash Table

key	record	hash Table
10	25	10
11	33	11
12	54	12

old record

old hash

new record

new hash

Bucket

$$33 \bmod 5 = 3$$

$$25 \bmod 5 = 0$$

$$54 \bmod 5 = 4$$

Basic Concepts in Hashing:

- ① Hash Table: Hash table is a data structure used for storing & retrieving data quickly. Every entry in the hash table is made using hash table function

2) Hash functions : -

- Hash funcⁿ is a funⁿ used to place data in hash table & to retr^o data from hash table.
- These hash funcⁿs are used to implement hash table.
e.g.-

3) Bucket: The hash function $h(key)$ is used to map several dictionary entries in the hash table. Each position of the hash table is called bucket.

4) Collision: Collision is a situation in which hash function returns the same address for more than one record.

e.g. $85 \bmod 5 = 0$

$85 \bmod 5 = 0$

- ⇒ **Probe:** Each ~~will~~^{has} calculateⁿ of address
→ rest for success is known
as a probe.
- ⇒ **Synonyms:** The ~~set~~^{of} keys that has ~~for~~^{to} the same location are synonyms. e.g. 25 & 55.
- ⇒ **Overflow:** When Hash table becomes full & no record needs to be inserted then it is called overflow.

Hash Functions :

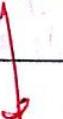
factor multiplied with with updating

Types

Linear

Hash funⁿ

long division, digit, logarithm



Division
method

Multiplication
method

Extracting
method

Raised
square
method

Folding
universe
method

Basic Concept Continue →

Overflow: the result of many keys hashing to a single address & lack of room in the bucket is known as an overflow. Collision & overflow are synonymous when the bucket size is fixed.

"Open or external hashing": When we allow records to be stored in potentially unlimited space, it is called as open or external hashing.

"Closed or internal hashing": When we use fixed space for storage eventually limiting the numbers of records to be stored, it is called as closed or internal hashing.

Hash fun": Hash fun" is an arithmetic "fun" that transforms a key into an address which is used for storing & retrieving a record.

Perfect hash function: The hash function that transforms different keys into different addresses is called a perfect hash function.

The worth of a hash function depends on how well it avoids collision.

Load density: The maximum storage capacity, that is the maximum no. of records that can be accommodated is called as loading density.

Load factor: Load factor is the number of records stored in a table divided by the maximum capacity of the table, expressed in terms of percentage.

1) Division method :

$$551 \div 10 = 4$$

$$72, 7 \div 10 = 2$$

$$89, 1 \div 10 = 9$$

$$37, 7 \div 10 = 7$$

Hash Table

0	
1	
2	
3	37
4	59
5	
6	
7	37
8	
9	89

2) Multiplication "Hash fun":

1) multiply the key 'k' by a const A where A is in range 0 to 1
Then extract the fractional part of ka .

2) multiply this fractional part by n & take the floor.

The above steps can be formulae as

$$h(k) = \lceil m \{ ka \} \rceil$$

fractional part.

Donald Knuth suggested to use
 $A = 0.10 \quad 0.61803398987$

eg - sort & balance with 61
 using method of 107, by assuming mod 50.

$$h(k) = [m + (107 + 0.6 \dots) k]$$

so that a balanced tree is suggested
 with modulus 50, mod 0.12 gives

0.12

$$\begin{aligned} h(12) &= 4.80 + 2 \cdot 0.12 \\ &= 5 \end{aligned}$$

That means that 107 will be placed
 at index 6th in hash table.

Implementation of insertion in hash table - C++
 To insert some elements in hash table
 we need to define all kinds of nodes
 to contain some library files
 like memory allocation header file
 with main file as main.h
 classes to work with in s i p

Extraction

In this method some digits are extracted from the key to form the address location in table.

e.g. Suppose first, third & fourth digit from left is selected for hash key.

(4) 9 (7) (8) 2 4

↓ ↓ ↓
4 7 8

→ at 478 location in hash table of size 1000 the key can be stored.

4) mid square

This method works in following steps.

- 1) Square the key
 - 2) Extract middle part of the result
This will indicate the location of the key element in the hash tab!
- Cond: Mid-square is used when the key size is less than or equal to 4 digits because if key is large then it is very difficult to store its square as it should exceed the storage limit.

Let key = CB11D

$$\begin{array}{c} \text{CB11D} \\ \downarrow \quad \downarrow \\ 9 \boxed{8} \boxed{7} 8 3 2 \end{array}$$

For the hash table of size of 1000
 $H(CB11D) = 783$

5) Folding

There are two folding techniques

i) fold - shift ii) Fold boundary.

i) fold - shift : In this method the key is divided into separate parts whose size matches with the size of required address. Then left & rightmost parts are shifted & added with middle part.

Example : Let key = CB11D

ii) Fold boundary : In this method the key is divided into separate parts. The leftmost & rightmost parts are folded on a fixed boundary & added with middle part.

Example : Let key = CB11D

Example : Let key = CB11D

key.

8 4 5 6 7 8 1 2 3

3 4 5 6 7 8 1 2 3



$$\begin{array}{r} 345 \\ + 678 \\ \hline 1123 \end{array}$$

Discard

3 4 5 6 7 8 1 2 3

digit
reversed



5 4 3

8 7 6

2, 3, 2, 1

1, 5, 4, 2, 1, 0, 9

Discard.

Index = 542

Index = 146.

b

* Position fold shifting. b) Fold Boundary.

* Properties of Good Hash Function

1. The hash function should be simple to compute.

2. No or few collisions should be likely while placing the lot of records in the hash table. Ideally, no collision should occur. Such a function is called perfect.

3. Hash funⁿ should produce such keys which will get distributed uniformly over array.
4. The hash funⁿ should depend on every bit of the key. Thus the hash funⁿ in that simply extracts the portion of a key is not suitable.

6. Rotation :

When the key is sent, they vary only in the last digit & this leads to the creation of synonyms. Rotation of the key would minimize this problem. Here, the key is rotated right by 1 digit & then folding technique is used to avoid synonyms.
e.g. 12060 When it is rotated we get 512060.

Universal Hashing:

The main idea behind universal hashing is to select the hash function at random at runtime from a carefully designed set of functions.

Because of randomization, the algorithm can behave differently upon each execution, even for the same input. This approach guarantees good average-case performance, no matter what keys are provided as input.

* Collision Resolution Strategies:

Three methods:

1. Open addressing (closed Hashing)

- a) Linear probing \rightarrow chaining
- b) Quadratic probing
- c) Double hashing
- d) Rehashing.

2. Separate chaining (Linked List.)

3. Bucket Hashing

inserted
same table. { open Addressing
closed Hashing
external Hashing

closed Addressing
open Hashing
internal Hashing

Q) Open Addressing:

This is collision handling technique in which the entire hash table is searched in systematic way for empty cell to insert new item if collision occurs.

Q) Linear probing:

When collision occurs i.e. when 2 records demand space for the same location in the hash table, then the collision can be solved by placing second record linearly down wherever the empty location is found.

$(key + i) \bmod \text{Max Table size}$

$$(2+1) \rightarrow 10$$

$$10 \bmod 4 = 2 \quad (\text{mod } 4)$$

e.g. Table of size 100 & key 1044
new key 2¹⁰ = 3544

key 1 is stored at location

$$\text{Hash}(key1) = \text{key1} \% 100$$

$$= 1044 \% 100$$

$$= 44$$

Index	key
0	
1	
2	
3	
4	1044
5	3549
6	
7	

Now have to store 2nd key = 3544

then Hash(key2) = 3544 mod 100

all next slot need add = 44. (as 3544)

But slot location 24 is already occupied so store it at next empty location i.e. 44.

so 3544 is stored at

location 24.

* [Problem with linear probing];

One major problem with linear probing is major primary clustering. Primary clustering is a process in which a block of data is formed in the hash table when collision is resolved.

eg.	0	39	cluster.
	1	29	
	2	9	
	3		$19 \cdot 1.10 = 9$
	4		$18 \cdot 1.10 = 8$
	5		$39 \cdot 1.10 = 9$
	6		$29 \cdot 1.10 = 9$
	7		$88 \cdot 1.10 = 8$
	8	18	
	9	19	

This

Linear probing using 2 ways.

i) With Replacement

ii) Without replacement

iii) Linear probing w/ with Replacement

6 20

eg. (24, 13, 16, 15, 19, 20,

22, 14, 17, 26.) 84 (ges)

2 22

3 13

4 24

5 15

6 16

Place 14 Here, cause $14 \cdot 1.10 = 4$

7 14

The index 4 contain the element

8 19

24. This is a proper record at its place.

9

Hence to place 19 we have to move down in search of an empty slot.

0	20	
1	26	← 26 is placed over here as empty space is available.
2	22	
3	13	
4	24	
5	15	
6	16	
7	17	
8	14	
9	19	

Next element $17 \cdot 1.10 = 7$ this position was occupied by 14. But 14 is not the proper record at this place. Hence we replace 14 by 17 & then place 19 at next empty slot.

Next element is 26, but $26 \cdot 1.10 = 6$ & 16 is a proper record at this place.

Hence at next empty location 26 is placed.

ii) Linear probing without replacement.

0	20	for linear probing without replacement
1	26	slot becomes full
2	22	insertion continues
3	13	insertion continues
4	24	full slot, move to next
5	15	insertion continues
6	18	
7	14	Collision occurs at index 7
8	17	4. Hence probing 14 at the next empty slot
9	19	and inserting 14 at index 8
		↓ at index 7, the 14 is already placed. Hence at the next empty slot 17 is placed.

→ $26 \text{ } \gamma. 10 = 60$! The collision occurs. Hence 26 is placed at next empty slot.

2) Quadratic Probing:

Quadratic probing operates by taking the original hash value and adding successive values of an arbitrary quadratic polynomial to the starting value. This method uses following formula.

$$\text{Value to Insert} = \text{Hash}(key) + i^2 \mod m$$

Here m is a prime number or any prime number.

e.g. Hash table of size 10 & 87, 90, 55, 22, 11, 17, 49, 87.

Index	Value
0	100
1	110
2	22
3	
4	
5	55
6	
7	37
8	
9	

Now if we want to place it a collision will occur as $17 \cdot 10 = 7$ and bucket 7 has already an element 87. Hence we will apply quadratic probing to insert this record in the hash table.

$$\text{Hickey} = (\text{Hash key}) + i^2 \cdot m$$

where $i = 0, 1, 2, \dots$

$$i = 0, \quad (17+0)^2 \cdot 10 = 7 \rightarrow 10$$

$$i = 1, \quad (17+1^2) \cdot 10 = 8$$

(previous slot is 10)

so 17 will get placed at 8.

After insertion an update will be done.

Then decrementing which will be placed at old index 8.

Now need to replace the 87.

Condition for replacement is 87

$$\text{if } (87+i^2) \cdot 10 = 7 \rightarrow \text{No}$$

$(87+1^2) \cdot 10 = 8$... but already occupied.

$(87+2^2) \cdot 10 = 17$... already occupied

$$(87+3^2) \cdot 10 = 26 \dots \text{This slot is free.}$$

process of insertion ends here.

No doubt some slots remain empty.

but they will not be all.

0	80
1	11
2	22
3	33
4	44
5	55
6	87
7	37
8	17
9	249
10	

3) Double Hashing

- Double Hashing is technique in which a second hash funⁿ is applied to the key when a collision occurs. By applying the second hash funⁿ we will get the number of positions from the point of collision to insert.

2 rules. 1. 2nd funⁿ.

- i) It must never evaluate to zero.
- ii) must make sure that all cells can be probed.

The formula \rightarrow $H_1(\text{key}) \rightarrow \text{hash}$

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

$$H_2(\text{key}) = m - (\text{key} \% m)$$

$$\text{Hash}(\text{key}) = H_1(\text{key}) + H_2(\text{key}) \quad \{ \text{not necessary.} \}$$

where m is a prime number

smaller than the size of the table.

e.g. Consider the elements

37, go, 45, 22, 97, 49, 53

$$H_1(37) = 37 \% 10 = 7$$

$$H_1(go) = go \% 10 = 0$$

$$H_1(45) = 45 \% 10 = 5$$

$$H_1(22) = 22 \% 10 = 2$$

$$H_1(97) = 97 \% 10 = 7$$

$$H_1(49) = 49 \% 10 = 9$$

$$H_1(53) = 53 \% 10 = 3$$

$$H_2(37) = m - (37 \% m)$$

$$H_2(go) = m - (go \% m)$$

$$H_2(45) = m - (45 \% m)$$

$$H_2(22) = m - (22 \% m)$$

$$H_2(97) = m - (97 \% m)$$

$$H_2(49) = m - (49 \% m)$$

$$H_2(53) = m - (53 \% m)$$

$$H_2(17) = m - (17 \% m)$$

Now if 17 is to be inserted

then

$$H_1(17) = 17 \% 10 = 7$$

$$H_2(17) = m - (17 \% m)$$

Here m is a prime number smaller than the size of the hash table, prime numbers smaller than m table size is 10 is 7.

Hence $m = 7$. Now we have to insert 37. That means we have to insert the element at 4 places from 37. In short we have to take 4 jumps.

\therefore the 17 will be placed at index 1.

How to insert 55.

$$H(55) = 55 \% 10 = 5 \quad (\text{Collision})$$
$$H(55) = 7 - (55 \% 7) = 7 - 6 = 1$$

That means we have to take one jump from index 5 to 55. Finally the hash table will be

0	90
1	17
2	22
3	
4	
5	45
6	55
7	37
8	15
9	49

4. Rehashing | Hash Table Overflow

- It is a technique in which the table is resized, i.e. the size of the table is doubled by creating a new table. It is preferable if the total size of table is a prime number.
- These are situations in which the rehashing is required.-
 - When table is completely full.
 - With quadratic probing when the table is filled half.

- When insertion fail due to overflow.

In such situations, we have to transfer entries from old table to the new table by recomputing the suitable hash for positions.

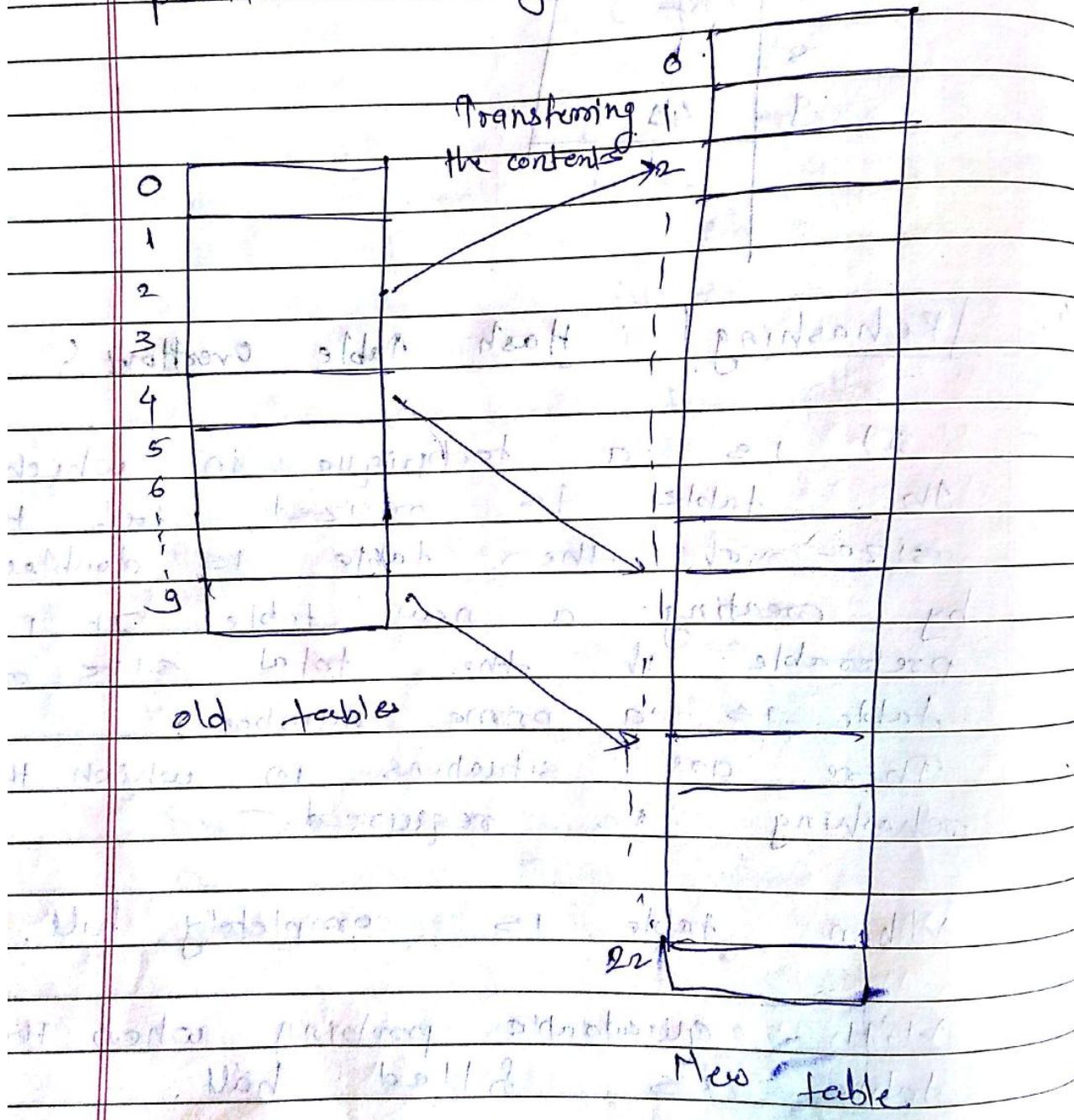


Fig : Rebasing.

e.g. $37, 90, 53, 22, 17, 49, 87$

$$H(key) = \text{key} \% \text{tablesize}$$

$$37 \% 10 = 7$$

$$90 \% 10 = 0$$

$$53 \% 10 = 3$$

$$22 \% 10 = 2$$

$$17 \% 10 = 7 \rightarrow \text{collision}$$

solved by linear
probing.

$$49 \% 10 = 9$$

$$87 \% 10 = 7$$

0	6	90
1	3	7
2	2	2
3	5	R
4	9	0
5	5	5
6		
7		37
8		17
9		49

Now this table is almost full &
if we try to insert more & more
elements then overflow will occur.

So will rehash by doubling the
table size. So double is 20 but
it is not prime no. so next is
23.

Theo

$$37 \times 123 = 14$$

$$90 \times 123 = 21$$

$$55 \times 123 = 9$$

$$22 \times 123 = 22$$

$$17 \times 123 = 17$$

$$49 \times 123 = 3$$

$$87 \times 123 = 18$$

0

1

2

49

1

9

55

1

22

1

17

1

3

1

2

1

8

1

7

1

4

1

2

1

9

0

1

2

2

8 digit numbers are added. If sum is more than 8 digits, then it is discarded.

Given 3 numbers. Total of 17. 17 is less than 8.

Given 3 numbers. Total of 18. 18 is less than 8.

Given 3 numbers. Total of 19. 19 is less than 8.

Given 3 numbers. Total of 20. 20 is less than 8.

Given 3 numbers. Total of 21. 21 is less than 8.

Given 3 numbers. Total of 22. 22 is less than 8.

W Examples

- 11 Give the input (4371, 1323, 6173, 4199, 4344, 9679, 1889) & hash function $h(x) = x \bmod 10$, show the

results for the following:

- i) Open addressing hash table using linear probing
- ii) Open addressing hash table using quadratic probing.
- iii) Open addressing hash table with 2nd hash fun $h_2(x) = 7 - (x \bmod 7)$

II program for i) Linear probing

II hash fun to get position.

```
int hashC(int Key)
```

```
return (key * Max);
```

3

II fun to insert new record.

```
int linear_probe(int hashtable[], int key)
```

2

```
int pos, i;
```

```
pos = hash(key);
```

if C Hashtable [pos] == 0) { // empty
}

 Hashtable [pos] = key;

 return pos;

}

else

{

 for (i = pos + 1; i < max; i++) { // pos+1 = pos; itt
 if C Hashtable [i] == 0)

{

 Hashtable [i] = key;

 return i;

}

 return -1; // overflow

}

Struct node

{

 int data;

 node *next;

}; h[10];

2] Chaining [Linked List]

In this technique, we can store the linked list inside the hash table, in the unused slot. This is used to handle synonyms using chaining; it contains together all the records that hash to the same address. Instead of relocating synonyms, a linked list of synonyms is created whose head is home address of synonyms.

As involvement of pointers is there some extra memory is needed for storing pointers.

0	100	
1		

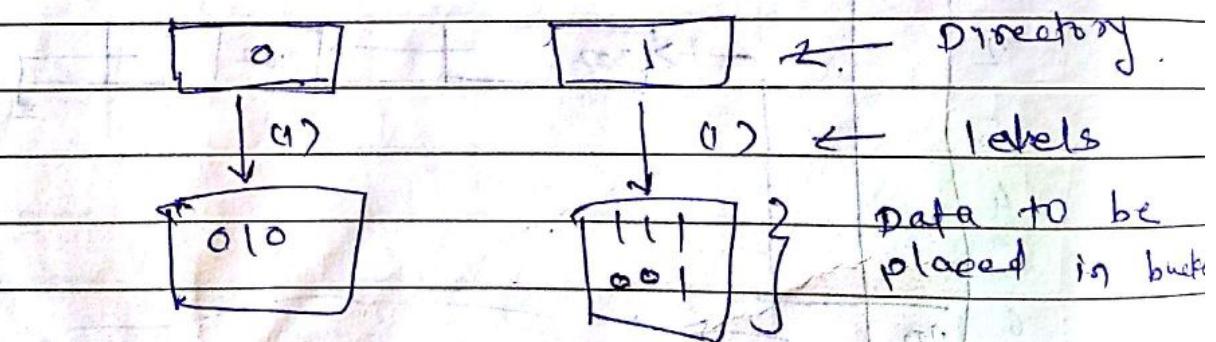
addresses.

Extendible Hashing:

- Extendible hashing is a technique which handles a large amount of data. The data to be placed in the hash table is split by extracting certain number of bits.

In extendible hashing referring to the size of directory the elements are to be placed in buckets. The levels are indicated in parenthesis.

e.g.



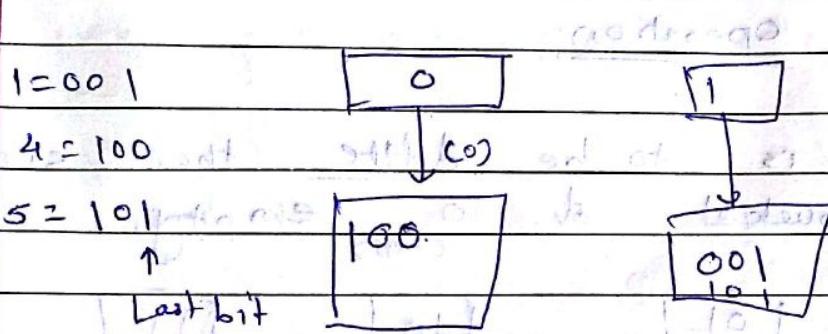
- The bucket can hold the data of its global depth.

If data in a bucket is more than global depth then, split the bucket and double the directory.

Insertion operation:

e.g. to insert 1, 4, 5, 7, 8, 10. Assume each page can hold 2 data entries (2 is the depth).

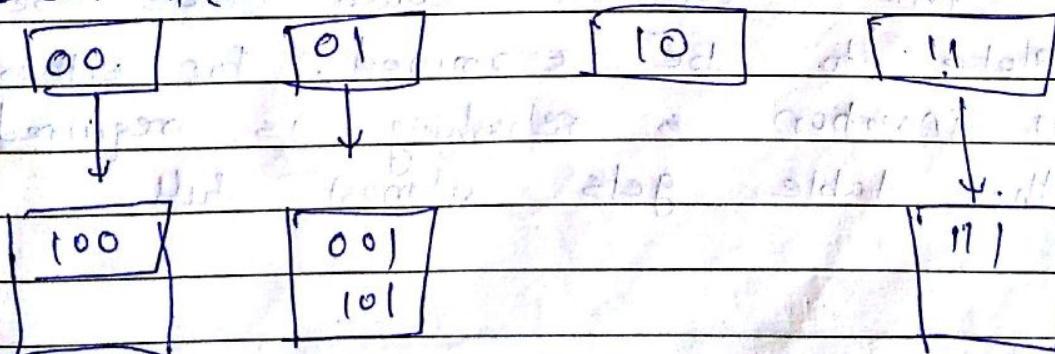
Step 1: Insert 1, 4.



We will examine last bit of data & insert the data in bucket.

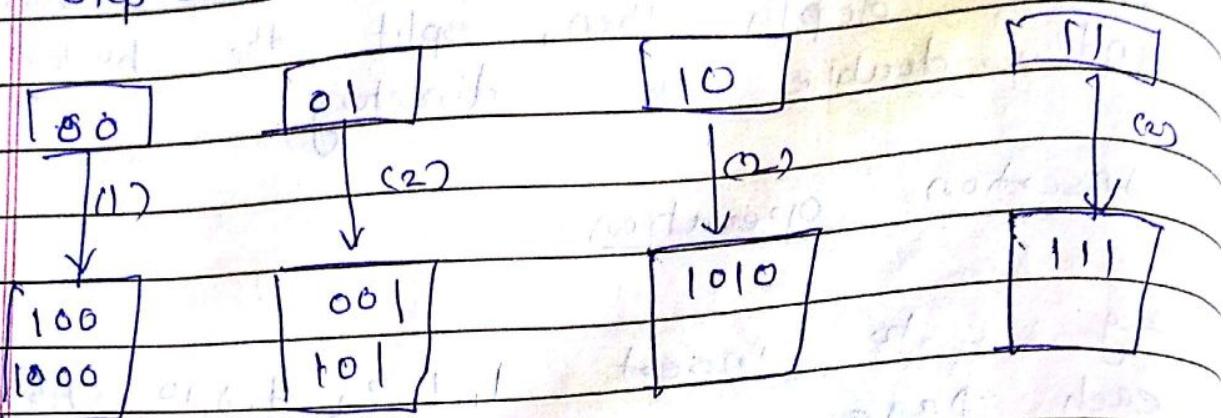
The bucket is full. Hence double the directory.

Insert 4 in this resulting tree.



Step 3 & 4:

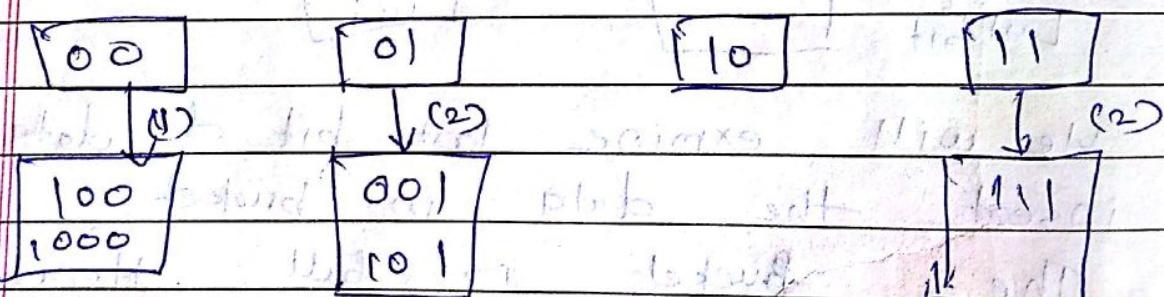
Insert 8 & 10. i.e. 1000 &



Thus data is inserted using extensible hashing.

Deletion Operation:

If 10 is to be deleted, then simply make bucket of 10 empty.



- The problem with open & closed Hashing is that collision could cause several slots to be examined. For either insertion or rehashing is required when the table gets almost full.

Closed Addressing:

[Chaining without Replacement]

e.g. 131, 31, 4, 21, 61, 6, 31, 8, 9.

replacement is not allowed

Index

key char

0 -1 start-not pushed

1 131 2

2 21 5 ! forward

3 3 3 -1 back at 3

4 4 -1

5 6 1 forward

6 6 6 forward down left

7 7 1 -1

8 8 -1

9 -1 forward

Chaining with Replacement

e.g. 131, 21, 31, 4, 52, 2, 9, 4

Index

key

char

0 -1

1 131 d 131

2 2+2 \$ -1

3 31 31 -1

4 4 4 -1

5 5 5 -1

6 21 21 3

7 21 21 -1

8 -1 -1

* Appl'n of Hashing

1. In Compilers to keep track of declared variables.
2. for Online spelling checking hashing functions are used.
3. Hashing helps in cache playing to store more pages.
4. For browser program while cache the web pages, hashing is used.

* Dictionary : It is a collection of key and value where every value is associated with the corresponding key.

* Load Density : It is denoted by $\frac{n}{b}$.

n = no. of keys

b = no. of buckets.

It should be less than $= 0.5$

$$n = 5, b = 10$$

$$\frac{5}{10} = 0.5$$

* Dictionary :
 is a collection of pairs of
 key and value where every value
 is associated with the corresponding
 key.

e.g. Library Database,

API for Dictionary.

instance : Dictionary is pair of key &
 value.

Operations :

size(), remove, contains, HA

Insert();

delete();

searches :

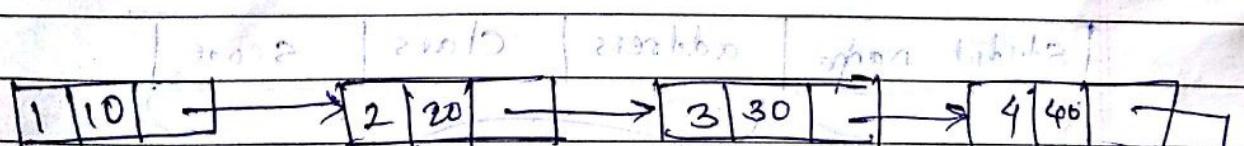
{ } : logo, registration, contacts, etc.

representations of dictionary

* Representation of Dictionary:

1) Sorted array (array sorted and sorted).

or Sorted chain (linked list).



Key ↑ value.

struct node

{
 int key;
 int value;
 node *next;

};

Difⁿ operatⁿ on Dictionary

1) Insert

2) Delete

3) Search

4) size.

All operations same as Linked

* Applⁿ of Dictionary:

1) student registration applⁿ:

A record of student can be deleted or updated from the dictio when a new student registers for course his/her record is inserted in the dictionary and counts feature

Student name	address	class	score
key	Value		

Identifiers - variable, funⁿ, arrays, struct.

PAGE NO.....

DATE.....

3) Telephone directory Dictionary

4) Word dictionary

5) Symbol Table

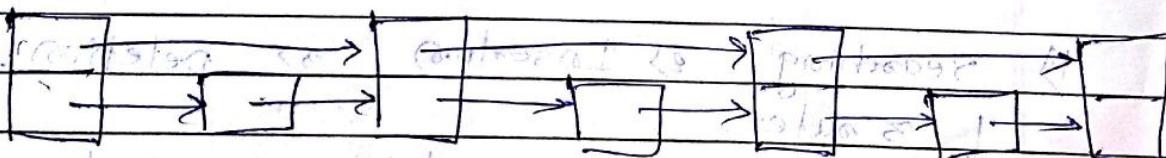
* Skip List

- Skip List is a variant list for the linked list.

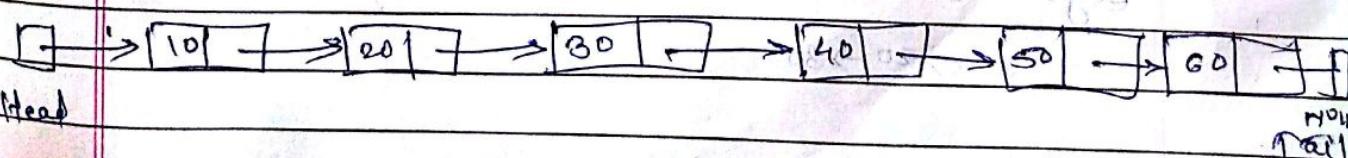
- Skip Lists are made up of a series of nodes connected one after the other.

The no. of references each node contains is determined randomly. The no. of references a node contains is called its node level.

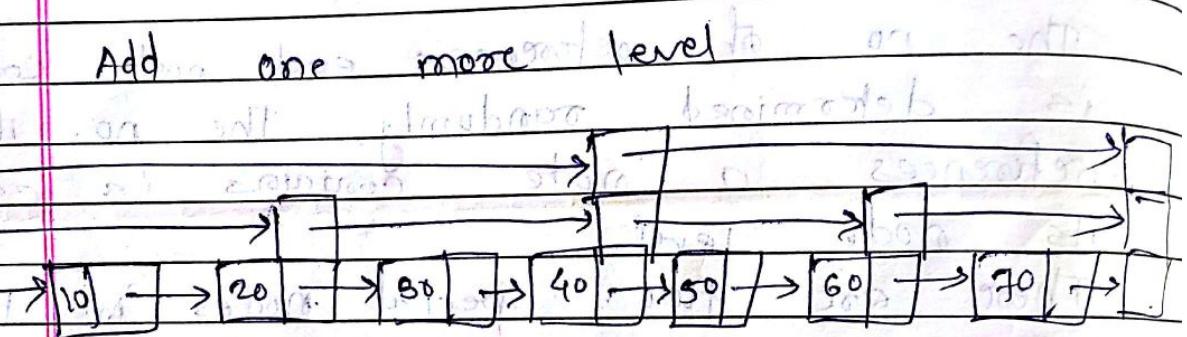
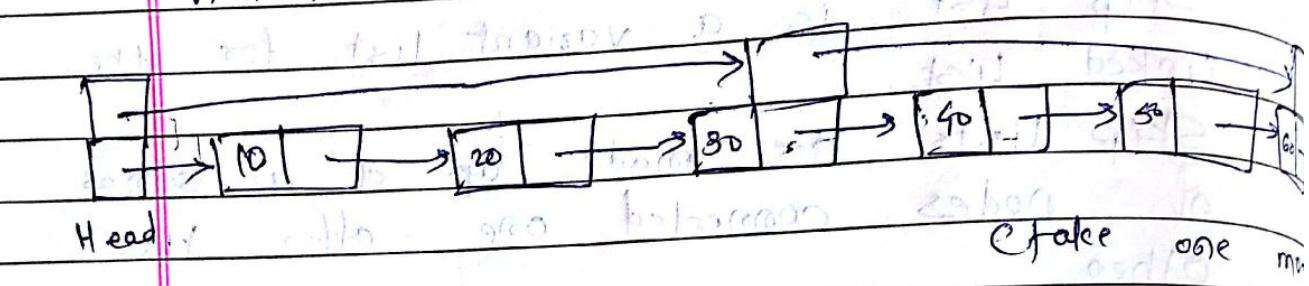
- There are two special nodes for the skip-list one is head node which is the starting node of the list & tail node is the last node of the list.



Consider a sorted chain as



Now to search any node from above sorted chain we have to search sorted chain from head node by each node. But this searching time can be reduced if we add level every alternate node. This extra level contains the forward pointer of node.



* Operations on skip list :

1) Searching, 2) Insertion, 3) Deletion.

- 1) If $\text{key} == \text{current} \rightarrow \text{data}$, then node is found.
 - 2) If $\text{key} < \text{current} \rightarrow \text{data}$, go down one level.
 - 3) If $\text{key} > \text{current} \rightarrow \text{data}$, go right along the same level.
- e.g. above e.g. search. 20

Hashing - Extra

Explain linear probing with & w/o replacement using the following data

12, 01, 04, 03, 07, 08, 10, 02, 05, 14, 06, 28.

Assume bucket from 0 to 8
each bucket has one slot.

Calculate average cost/no. of comparisons for both.

Solⁿ: We assume hash funⁿ = key mod 10. As buckets are from 0 to 8 & each bucket has one slot. the keys can be inserted in the hash table using our hash function as follows.

Index	key
-------	-----

0	10
---	----

1	01
---	----

2	12
---	----

3	03
---	----

4	04
---	----

5	02
---	----

6	05
---	----

7	07
---	----

8	08
---	----

How hash table is filled
can not insert 06
in the hash table

key	12	01	04	03	07	08	10	02	05	14	06	28
No. of	1	1	1	1	1	1	1	4	02	6	10	10

Comparisons.

Thus total no. of comparisons

$$= 1 + 1 + 1 + 1 + 1 + 1 + 0.2 + 0.2 + 6 + 10 + 10$$

$= 39$ buckets are examined

Load density $\alpha = \frac{n}{b}$

$$= \frac{12}{10} > 1.2$$

The avg. cost can be computed by calculating successful search (S_n) & unsuccessful search (U_n).

$$S_n = \frac{1}{2} \left(1 + \frac{1}{(1-\alpha)} \right)$$

$$\checkmark \frac{1}{2} \left(1 + \frac{1}{(1-1.2)} \right)$$

$$S_n = -2$$

$$U_n = \frac{1}{2} \left(1 + \frac{1}{c(1-\alpha)^2} \right)$$

$$= \frac{1}{2} \left(1 + \frac{1}{(1-1.2)^2} \right)$$

$$= \frac{1}{2} \left(1 + \frac{1}{0.04} \right)$$

$$= 1.8$$

$$\text{Total cost} = S_n + U_n$$

$$= -2 + 1.8$$

$$= 11$$

8