

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

Hashing

MAY

24

18th Week
125-241

04

SATURDAY

→ It is mainly used to implement dictionary, where you have key value pair, as well it also used when you set up a database.

08

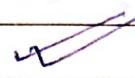
→ Search

09

Insert $\rightarrow O(1)$ on average

10

delete



→ Not useful for sorted data to be searched using binary search.

11

(a) finding closest value of AVL trees (or) Red black tree.

12

(b) sorted data (or) binary search (or) (c) prefix searching (or) (d) the sorted & unsorted data in

01

Application of Hashing

(1) Dictionaries

(2) Database Indexing

(3) Cryptography

(4) Caches

(5) Symbol tables in compilers/interpreters

(6) Routers

(7) Caches

(8) getting data from database need to sort it

and go on, follow these steps before hashing

02

2nd most used DS is Computer science.

03

Direct Address Table

DAT is a DS that has the capability of mapping record to their corresponding keys using arrays. In direct address tables, records are placed using their key values directly as indexes.

These facilitate fast searching, insertion and deletion operation.

MAY

05

'24

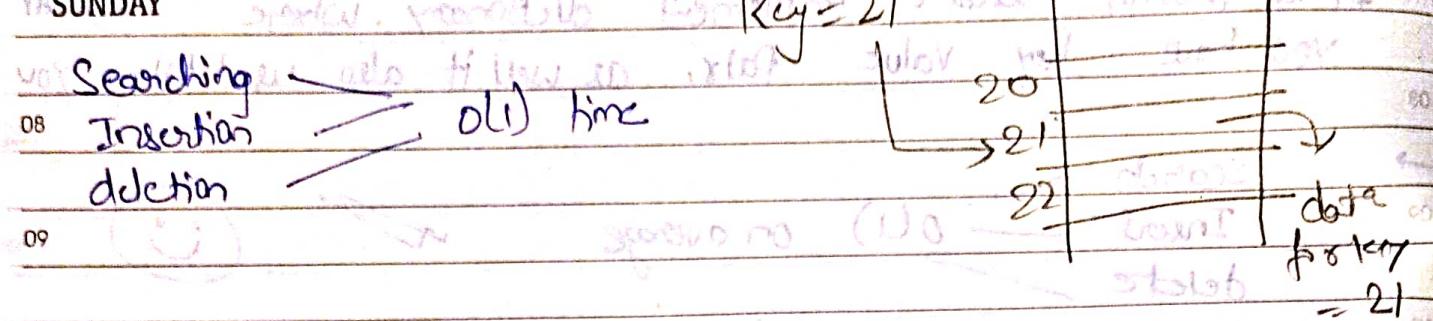
19th Week
126-240

SUNDAY

midterm

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

MAY 2024



10 Limitations :-

- (1) Prior knowledge of maximum key value.
- (2) Practically useful only if the maximum value is very less.
- (3) It causes wastage of memory spaces if there is a significant difference b/w actual records and maximum value.

01 Hashing can overcome these limitations of direct address tables.

Hash Functions

- Hashing is a technique (Or) process of mapping keys and values into the hash table by using a hash function. It is done for faster access to elements depend on the position of object factors.
- The efficiency of the hash function.

05 The task of hash function to convert those large Keys, String, floating point Number into small values that can be used as an index in the hash table.

- (a) Should always map a large key to some small keys.
- (b) Should generate values from 0 to m-1.
- (c) Should be fast, O(1) for integers and O(len) for string of length 'len'.
- (d) Should uniformly distribute large keys into hash tables slots.

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

MAY

'24

19th Week
127-239

06

MONDAY

Example Hash function.

- ① $h(\text{large_key}) = \text{large_key \% m}$
- ② for string weight sum $\text{str} = 'abcd'$

$$(str[0] * n^0 + str[1] * n^1 + str[2] * n^2 + str[3] * n^4) \% m$$
- ③ universal hashing.
 - ↳ We choose Prime No. which are not close to powers.
 - Set of hash function, we have to pick randomly and your hash function

Collision handling

Birthday Paradox
27 more

Collision is the conditions in which two keys map to the same entry in the hash table.

- ④ If we know keys in advance, then we can go with perfect hashing.
 - ⑤ If we don't know Keys, then we use one of the following.
- Following are the ways to handle collision
- chaining
 - open addressing
 - linear probing
 - Quadratic probing
 - double hashing

- Chaining (Separate chaining)
- The idea behind separate chaining is to implement the array as a linked list called a chain.
 - The linked list data structure is used to implement this technique.
 - what happens is: when multiple elements are hashed into the same slot index, then these elements are inserted into a singly linked list which is known as chain.

MAY

07

'24

19th Week
128-238

S	M	T	W	T	F	S
5	6	7	8	9	10	11
12	13	14	15	16	17	18
19	20	21	22	23	24	25
26	27	28	29	30	31	MAY 2024

TUESDAY

 $\text{hash}(\text{key}) = \text{key \% } 7$ $\text{keys} = \{50, 21, 58, 17, 15, 49, 56, 22, 23, 8, 5\}$

08

0	21	→	49	→	56
1	50	→	15	→	22
2	58	→	30	→	23
3	17	→	9	→	8
4	25				
5	49				
6					

12 Performance from last out, division operation with multiplication

 $m = \text{No. of slots in hash table}$ 01 $n = \text{No. of keys to be inserted}$ Load factor $\lambda = n/m$

02

Expected chain length = λ 03 Expected time to search = $O(1 + \lambda)$ Time for insert / delete = $O(1 + \lambda)$

04

Data Structure for storing chains

05 (not cache friendly) (a) linked list — search $O(l)$, Delete $O(l)$, Insert $O(l)$ 06 (not cache friendly) (b) Dynamic size arrays { vector in C++, Array list in Java } $T_c = O(1)$

07 (not cache friendly) (c) Self balancing of BST (AVL tree, Red Black tree)

 $O(\log l)$ — search

insert

delete

SUN	MON	TUE	WED	THU	FRI	SAT
30	1					
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29						

JUN 2024

MAY

124

19th Week
129-237

08

WEDNESDAY

Implementation of chaining in Python

BUCKET = 7

08

0	[]		70	[]	56	[]
1	[]		71	[]		[]
2	[]		9	[]	72	[]
3	[]	Insert →	[]	[]		[]
4	[]	70, 71, 9, 72	[]	[]		[]
5	[]	56	[]	[]		[]
6	[]		[]	[]	56	[]

12

[[],[],[],[],[],[],[]]

hash = n % BUCKET

h = MyHash(7)

class MyHash:

h.insert(70)

def __init__(self, b):

h.insert(71)

self.BUCKET = b

h.insert(9)

self.table = [[] for i in range(b)]

h.insert(56)

def insert(self, n):

h.insert(72)

i = n % self.BUCKET

print(h.search(56))

self.table[i].append(n)

h.remove(56)

def remove(self, n):

print(h.search(56))

i = n % self.BUCKET

o/p - True

self.table[i].remove(n)

False

def search(self, n):

i = n % self.BUCKET

return n in self.table[i]

06

07

MAY

'24

09

19th Week
130-236

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

THURSDAY

Frequency of Array Elements

08 i/p - arr E] = {10, 20, 12, 10, 15, 10, 12, 12} o/p - 10 3

09 i/p - 12 3

15 1

10 i/p - 20 1

11 i/p - arr E] = {10, 10, 10, 10} o/p - 10 4

i/p - arr E] = {10, 20} o/p - 10 1 20 1

① def countfreq (arr, n):
01 mp = dict ()
for i in range (n):02 if arr [i] in mp . keys ():
if arr [i] in mp [arr [i]] + = 103 else:
mp [arr [i]] = 104 for m in mp :
(m, mp [m])

05 print (m, mp [m])

06 T.C = O(n)

07 O(n) d.o. p.b

T.C = O(n)

[f] of or not in a set

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

JUN 2024

MAY

124

19th Week
131-235

10

FRIDAY

~~Set &~~

(S?)

08. A set is an unordered collection data type that is iterable, mutable and has no duplicate elements.

- 09. (a) Distinct elements
- (b) Unordered
- (c) No indexing
- (d) Union, intersection, set difference etc are fast.
- 11. (e) Uses hashing Internally

12. S.add() \rightarrow S = {2, 4, 6, 8}, S2 = {3, 6, 9} \rightarrow S = {2, 3, 4, 6, 8, 9}

S.update() \rightarrow S = {2, 3, 4, 6, 8, 9} \rightarrow print(S1 | S2) \rightarrow {2, 3, 4, 6, 8, 9}

01. S.discard(b) \rightarrow S = {2, 3, 4, 6, 8} \rightarrow print(S1 & S2) \rightarrow {6}

S.remove(b) \rightarrow S = {2, 3, 4, 6, 8} \rightarrow print(S1 - S2) \rightarrow {2, 3, 4}

02. S.clear() \rightarrow S = {} \rightarrow print(S1 ^ S2) \rightarrow {2, 3, 4, 6, 8, 9}

len(S)

03. Intersection and difference function \rightarrow Symmetric difference function \rightarrow S1.difference(S2)

Union

 \hookrightarrow S1.union(S2)

S1.intersection(S2)

S1.symmetric_difference(S2)

04. Dictionary is a collection of key: value pairs \rightarrow Key: Value pair

05. Dictionary is a collection of key: value pairs \rightarrow Key: Value pair

06. Dictionary is a collection of key: value pairs, used to store data value like a map, which unlike other data types which holds only single value as an element.

- (a) Unordered
- (b) All keys must be distinct
- (c) Value may be repeated
- (d) uses hashing Internally

d = {10: 'xyz', 101: 'abc', 105: 'bcd', 104: 'abc'}

\rightarrow Similar to hashmap in Java, unordered map in C++.

MAY

11

'24

19th Week
132-234

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

SATURDAY

(C.R.)

(start)

- 08 d.get(101) → my2
- 09 d.get(125) → None
- 10 len(d) → 4
- 11 d.pop(105) → {110: 'abc', 101: 'wing', 106: 'bcd'}
- 12 dcl d[106] →
- 13 d.popitem() → pairabt 011
- 14 ~~Find over the search tag, and then, which (b)~~
- 15 ~~Open addressing~~

12 Like separate chaining, open addressing is a method for handling collision. 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. (Note that we can increase table size by copying old data if needed). This approach is also known as closed hashing.

03 The entire procedure is based upon probing, we will understand the types of probing ahead. (Cache friendly)

(2) a) Linear Probing.

05 Linearly search for next empty slot.

$$\text{hash}(\text{key}) = (\text{int}) \ \% \ \text{table_size};$$

06 ~~Search slot, start of searching in circular manner for next empty slot~~

0	49
1	50
2	51
3	16
4	56
5	15
6	19

$$\text{hash}(\text{key}) = \text{Key} \% 7 \quad \text{00 to 10}$$

07 Insert(50), insert(51), insert(15),
Search(15), Search(64), delete(15)
Search(15) printout 820

08 ~~Search(50), search(51), search(15), search(64), delete(15)~~

09 ~~insert(50), insert(51), insert(15), search(15), printout 820~~

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

JUN 2024

MAY

'24

20th Week
133-233

12

Search — We compute hash function we go to that index and compare if we find, we return true; otherwise (d)

08 → we linearly search.

→ We stop when one of the three cases arise

09 (i) empty slot → False

10 (ii) "key" found → True

11 (iii) Traversed through the whole table → False

11 Delete — Problem with simply making slot empty when we delete.

→ So we don't mark the slot empty we mark it as

12 now officially deleted or tombstone slot or silent slot

→ (Tombstone slot) = (Dead slot) = Local destruction

02 → Primary clustering vs problem with linear probing.

→ If all slots are filled then it fails.

02 To handle primary clustering problem with linear probing is

hash(key) = key % 0.7 with random initialization &

03 but result hash(key, i) = (hash(key) + i) % 0.7 + 1

04 (1) Quadratic probing → (i) $h_1(\text{key}, i) = (\text{hash}(\text{key}) + i^2) \% m$

→ (ii) Double hashing → $h_1(\text{key}, i) = (\text{hash}_1(\text{key}) + i \times \text{hash}_2(\text{key})) \% m$

05 (1) $\alpha < 0.5$ if load factor is less than 0.5 then only

06 you will able to find empty slot.

(2) m is prime in (n/m) .

if two these things happen together then quadratic probing

works and they will finds empty slot.

MAY

13

'24

20th Week
134-232

S	M	T	W	T	F	S
5	6	7	8	9	10	11
12	13	14	15	16	17	18
19	20	21	22	23	24	25
26	27	28	29	30	31	MAY 2024

MONDAY

Topic: Hashing (Open Addressing) - Quadratic Probing

(b) Quadratic Probing: An hash with following form

- Quadratic Probing is a method with the help of which we can solve the problem of clustering that was discussed above.
- This method is also known as the mid square method.
- In this method, we look for the i^2 th slot in the i^2 th iteration. We always start from the original hash location. If only the location is occupied then we check the other slots.

(c) Double hashing: An hash with following form

- Double hashing is a technique that reduces clustering in an optimized way. $\text{hash}(\text{key}, i) = (\text{hi}(\text{key}) + i \cdot \text{h}_2(\text{key})) \% m$
- In this technique, the increments for the probing sequence are computed by using another hash function.
- We use another hash function $\text{h}_2(m)$ and look for the $i \times \text{h}_2(m)$ slot in the i^2 th iteration.
- If $\text{h}_2(\text{key})$ is relatively prime to m , then it always finds a free slot if there is one.
- (b) Distributes keys more uniformly than linear probing and quadratic hashing.
- (c) No clustering.

$$\text{hash}(\text{key}, i) = (\text{hi}(\text{key}) + i \cdot \text{h}_2(\text{key})) \% m$$

49, 63, 56, 52, 154, 148, 48

0 49

 $m = 7$

1

2 54

3 63

4 56

5 52

6 48

$$56 \rightarrow (0 + (x+1)) \% 7$$

$$52 \rightarrow (3 + (6 - (52 \% 6)))$$

$$(3 - 2)$$

obtained 3 slots out of 7 slots available

1 ✓

54 $\rightarrow 5 + (6 - 54 \% 6)$

$$5 + 6 = 11 \% 7 = 4$$

SUN	MON	TUE	WED	THU	FRI	SAT
30					1	
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29						

MAY

'24

20th Week
135-231

14

TUESDAY

- # Performance Analysis of search
- ⑥ Cache performance of chaining is not good because when we traverse a linked list, we are basically jumping from one node to another, all across the computer's memory.
- like chaining, the performance of hashing can be evaluated under the assumption that each key is equally likely to be hashed to any slot of the table.

$$m = \text{No. of slots in hash table}$$

$$n = \text{No. of keys to be inserted in hash table}$$

$$\text{Load factor } d = \frac{n}{m}$$

Effected time to search / insert / delete ≤ 1

$$d = \frac{1}{1 - (1 - \alpha)}$$

$\alpha = \text{fraction of the table left empty.}$

$$\alpha = 1 - d$$

Unsuccessful search means when you have either traversed through the whole table or you ended up with an empty slot.

Chaining

- | | |
|---|--|
| ① It is simple to implement | ② It requires more computation. |
| ③ Hash table never fills up; we can always add more elements. | ④ It may become full. |
| ⑤ Less sensitive to hash function | ⑥ Extra care required for clustering. |
| ⑦ Cache performance of chaining is not good as keys are stored using linked list. | ⑧ It provides better cache performance as everything is stored in same table. |
| ⑨ Extra spaces are required for linking ($L + \alpha$) | ⑩ might be needed to achieve same performance $= [1/(1-\alpha)]$ |
| → better technique | → old technique |
| → It is used when frequency and No. of keys are known. | When poverty comes in the door, love goes out the window. If keys are known, frequency may dynamically increase. |

MAY

15

'24

20th Week
136-230

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

MAY 2024

WEDNESDAY

Implementation of open addressing in Python

```

08 class HashTable:
09     def __init__(self, cap=10):
10         self.size = cap
11         self.table = [-1]*cap
12
13     def hash(self, n):
14         return n % self.size
15
16     def search(self, n):
17         if n == -2:
18             return False
19         h = self.hash(n)
20         t = self.table
21         i = h
22         while t[i] != -1:
23             if t[i] == n:
24                 return True
25             i = (i+1) % self.size
26         return False
27
28     def remove(self, n):
29         if n == -2:
30             return
31         h = self.hash(n)
32         t = self.table
33         i = h
34         while t[i] != -1:
35             if t[i] == n:
36                 t[i] = -2
37             i = (i+1) % self.size
38
39     def insert(self, n):
40         if n == -2:
41             return
42         if self.size == self.cap:
43             print("Overflow")
44             return
45         if self.search(n) == True:
46             print("Element already exists")
47             return
48         h = self.hash(n)
49         t = self.table
50         i = h
51         while t[i] not in [-1, -2]:
52             i = (i+1) % self.size
53         if t[i] == -1:
54             t[i] = n
55             print("Element inserted")
56             print(h)
57             print(h.search(56))
58             print(h.remove(56))
59
60             Whatever you are able to secure from a burning house is again
61             whatever you need to
62             return True
63             print(h.search(56))
64             print(h.remove(56))

```