

for hashing, the order mentioned is on avg., & not the worst case. for worst case these orders can blow upto  $n$ .

Date: / /  
Page No.:

## Dynamic Data Set (DDS)

We want to perform the following operations:-

→ Add ( $x$ )       $O(1)$        $O(1)$        $O(n)$        $O(1)$        $O(\log n)$

→ Delete ( $x$ )       $O(1)$        $O(1)$        $O(n)$        $O(1)$        $O(1)$        $O(\log n)$   
 $O(n)$        $O(n)$

→ Search ( $x$ )       $O(n)$        $O(n)$        $O(\log n)$        $O(1)$        $O(\log n)$

linked list (LL)	Array (A)	Sorted array (SA)	Hashing	Balanced Binary Search Tree (BBST)
			$O(1)$ is avg not the worst case.	$O(1)$ is avg not the worst case.

↳ doesn't help much

∴ to delete  $X$ , you have to 1st search for  $X$  & moreover in array you have shift everything by one position, after deleting.

needed only for sorted array

for general array delete  $X$  & put the last element at the position of  $X$ , since the order doesn't matter

Ex:-  $\boxed{6} \boxed{8} \boxed{9} \boxed{2} \boxed{1}$  & we have to delete '8'

the final array will be  $\boxed{6} \boxed{1} \boxed{9} \boxed{2}$

Order  $n$  is really bad, if we do  $n$  such operations, it suddenly blows upto order  $n^2$ .



In BBST, all the 3 operations can be done in order log n, but if you just want to do add, delete & search, it's still worth taking risk & dealing with Hashing; probably it's NOT worth going for BBST.

But there are more operations that we want to do, we will see more examples later on. However one of the operations worth mentioning is Rank(x)

Rank(x) will return the no. of numbers  $\geq X$  for a dynamic data set

We keep adding, deleting & there is also a search.

Hashing	BBST
Rank(x)	$O(n)$

To start with Hashing let's look at a concrete example. We're given a set of telephone numbers



Let's say that no. of telephone numbers =  $10^6$

10 digit

This is not too large a no. b/c the no. of airtel users in Bang. is much larger than this,

So we want to maintain a data str. which is dynamic, you should be able to add new no., delete any existing no. if needed & return true or false if asked that a certain no. exists in tel. directory or not.

First, now look at a similar ex. we have a set of internal (college) telephone no. Here the no. of users  $\leq 10^3$ . Hence, a three digit no. will do. ~~the work. However, let's say that the no. of user  $\approx 200$~~

Now, here we want to:-  
add, delete & search

To do this we can use an array of size 1000 (say). If a no. corresponding to an index exists then set it to true, otherwise false.

ex:- 268, 128, 438 are the no.

0	1	2	128	268	438
F	F	F	T	T	T

$\therefore$  the amount of space needed = 1000

Add. ] each of them happens in  
Delete ] order 1.  
search ]

It's just a table lookup.

However, if we compare this method for the prev. ex., the no. of telep. no. ( $10^6$ ) is not a problem. We can declare a array of size  $10^6$ . But if we want to apply the same method as mentioned above we'll need an array of size  $10^{10}$ . (NOT possible)

Which is where hashing comes into picture.

Formal definition of hashing:-

hashing is a "fun" from a set  $S$  (a dynamic set),  
 {for all practical purposes we'll assume  $S$  as an integer set. Ideally this is an infinite set.} to a set  $A$ .

Q. Why can we assume that any data is an integer?

Even if we're given a book, we want to interpret that lets say as an integer. Even if you give me a 2000 pg book the way computer stores it is in the binary representation.

$$A = \{0, 1, \dots, m-1\}$$

$$h: S \rightarrow \{0, 1, \dots, m-1\}$$

where  $m$  is the size of the table.

Ex? ;  ~~$h(x)$~~

Examples of hash fn:-

$$i) h(x) = (ax + b) \bmod m$$

$1 \leq a, b \leq m$

In particular you can take

$$b=0 \& a=1$$

$$\Rightarrow h(n) = n \bmod m$$

frequently used

This is a frequently used hash fn.

ii)  $h(x) = \lfloor m(ax \bmod 1) \rfloor ; 0 < a < 1$

a fractional no.  
bet 0 & 1

i.e.  $a$  is a fractional no.. It may be  
 $\frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \dots$

but this  
is NOT  
frequently  
used

$(ax \bmod 1)$  will give us the fractional part of  $a$ , leaving the integer part of it.

for ex:-  $2.87 \% 1 = .87$

Now, in the 10 digit teleph. no. ex.  $n$  i.e. the no. of tel. no. is  $10^6$ .

Q. The natural ques. arises is what should be the size of  $m$ ?

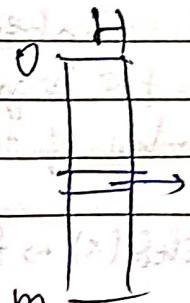
Sol<sup>n</sup> Typically, we take  $m = 2n$

Q. Why is it so?

Sol<sup>n</sup> We want to store  $n$  no., we create a space of size  $2n$ , so we have enough space to store these no. If we just take  $m$  smaller than  $n$ , then we can't obviously store everything in the table. So, just ensure that you make relatively a bigger table, so that we're able to store all those no.

Now, how

to add  $A(x)$ , we look at  $h(x)$   
& store  $x$  at  $H[h(x)]$   
i.e. at the index  $h(x)$



$\text{Delete}(x) \rightarrow H[h(x)] = -1$

i.e. go to the index  $h(x)$  & make it as  $-1$ ,  
 $-1$  indicates that there is nothing. ( $\because$  tel. no. are  
 always +ve, so -ve no. are not possible)

$\text{Search}(x) \rightarrow \text{check if } H[h(x)] = x$

If this is true then return true  
 otherwise, return false

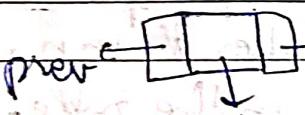
However, there exists a problem of collision when  
 we add like this.

collision is when  $x \neq y$  but  $h(x) = h(y)$

To handle the problem of collision, we've techniques  
 called collision resolution techniques:-

one of which being is the ~~de~~claring array as  
 an array of linked list (instead of simple array),  
 where every index is pointing to NULL when empty

structure of each node will look like



Now, when we compute  $h(x)$ , it will give a pointer  
 to a linked list. Adding  $x$  will then be basically to add  
 at the beginning of the linked list where index correspondingly  
 index of the array will be the head of the linked list

$\text{Delete}(x)$  is find  $h(x)$ , traverse through LL & delete that  
 node when you find it.

Search( $x$ ) is find  $h(x)$ , traverse through the LL & return true if you find the node  $\&$

Here, we have a problem that if a ~~naughty~~ <sup>naughty</sup> person knows the our hash fn<sup>n</sup> i.e.  $a, b$  then he can provide us with a data where each one has the same has value.

Real life ex:- About in 1995, Oracle had a database based on a hash fn<sup>n</sup> & everybody knew what the hash fn<sup>n</sup> is. So, on one fine day attackers hashed everything which mapped to the same location, the entire database crashed.

So, declare what your hash fn<sup>n</sup> is, is probably not the right way to do. So, if we have a large set of hash fn<sup>n</sup> then I can uniformly randomly pick a hash fn<sup>n</sup> & work with that.

So your program only says that I'm going to pick a random has fn<sup>n</sup> from a huge collection of hash fn<sup>n</sup>, but what that hash fn<sup>n</sup> is, nobody knows. That is decided on the runtime by the machine.

this collection is called a universal collection of hash fn<sup>n</sup>.

What is expected is that the no. of collision should be constant. Only then add, search & delete become order 1. Even if one of the location in hash table becomes  $\frac{1}{2}$  or  $\frac{1}{5}$  then the <sup>performing</sup> operations will become a linear time.

i.e. everytime we run the program a new hash table is created, <sup>from scratch</sup> so there should be a way to store it.

IF YOU WANT TO AVOID OPEN-ADDRESSING THEN YOU CAN USE EITHER OF THE FOLLOWING,



A.) Linear probing → Simplest of all

EITHER OF THE FOLLOWING,

3:-

Let  $A_1, A_2, A_3, \dots, A_8$  have the same hash value 5

then we will store them like this

5	
	$A_1$
	$A_2$
	$A_3$
	$A_7$
	$A_8$
	-1

Now, suppose we want to delete  $A_7$  then we will go to index 5 & search for  $A_7$  until we find it or we find (-1).

But this can also cause a problem.

Suppose  $A_6$  is deleted then the table will look

like this

& now we want to search for  $A_7$ .

At 11th index the search will stop since we get -1.

So, to handle this each index

should have <sup>also</sup> a flag value of True or False.

5	<del><math>A_1</math></del>
6	<del><math>A_1</math></del>
7	-1
8	<del><math>A_7</math></del>
9	<del><math>A_8</math></del>

If the flag is true then we will know that 5th is deleted from that ~~place~~ place and hence we continue our search.

Now if you want to search for a no.  $x$  then go to  $h(x)$  then go on searching for  $x$  until you find  $x$  or a  $(-1)$  with a false flag.

when you delete an  $x$  replace  $x$  with  $(-1)$  & Set the flag to true.

when you add, go to index  $h(x)$ , if it's occupied just keep going forward until you find  $(-1)$  (doesn't matter whether flag is true or false) & insert  $x$ . If, when moving forward, you reach the end of the array, come to ~~repeatedly~~  $0^{th}$  index. This is not a very efficient way to continue the procedure. Then do  $(\text{mod } m)$ .

This is not a very efficient way when you consider the search & other thing but as far as the space is considered, it will only take  $n$  space, you don't have to create LL & other things.

### B) ~~Linear Probing~~

#### Quadratic Probing

Everything is same like linear probing, except that instead of moving forward by one position we move forward by  $i+i^2$

Ex:-

5	A <sub>1</sub>
7	A <sub>2</sub>
11	A <sub>3</sub>

Initially  $i=0$  then  $i=1$

~~h(A<sub>1</sub>) + 1 + 1<sup>2</sup> = h(A<sub>1</sub>) + 3~~  
 $A_2$  will go to  $h(A_2) + 1 + 2^2 = h(A_2) + 9$   
 $A_3$  will go to  $h(A_3) + 2 + 3^2 = h(A_3) + 16$   
 Here also we've to take care of flow.

If  $h(n) + i + i^2$  goes out of range then  
do  $(\text{mod } m)$ .

### C.) Double hashing

Everything stays same as quadratic probing,  
except that instead of moving forward by  
 $i + i^2$ , we move forward by  $i h_1(x)$ , where  
 $h_1(x)$  is another hash fn.

If  $h(x) + i h_1(x)$  goes out of range then  
do  $(\text{mod } m)$ .

mind that  $h_1(x)$  can be -ve for any  $x$ ,  
instead of moving forward, actually we will  
backtrack. It's completely normal. Don't get  
psyched.

Even after knowing all these things, someone  
can give us a collection of no. which will map  
to the same location. In the open-addressing,  
it all becomes a LL. In the cases above, it all  
becomes an array. However, operations in both  
ways will take  $\Theta(n)$  time.

So, that's why we want to choose the hash fn  
uniformly at random. In order to do this, we  
should have a large collection of hash fn with  
some nice property. This collection is called  
universal hash fn.

## UNIVERSAL HASHING

$\mathcal{H}$  → set symbol for collection of hash fn<sup>n</sup>.

↓  
~~another set.~~

an infinite set of hash fn<sup>n</sup> which all map to  
 $\{0, 1, \dots, m-1\}$

This set is a universal hash fn if it satisfies the following property (s):—

$\forall x \neq y \quad h(x) = h(y), \quad h \in \mathcal{H}$

The no. of such hash fn  $\leq \frac{|\mathcal{H}|}{m}$   
 in  $\mathcal{H}$  for any pair  $x, y$ .

i.e. if the size of the hash table is smaller, then the no. of hash fn where the collision actually happens is more.

so, the no. of hash fn, where there will be a collision is directly proportional to the ~~no.~~ total no. of hash fn that you have & inversely proportional to the size of the hash table

of hash fn<sup>n</sup>

Q.

If we are given a set, then what is the probability that there will be a collision at a position in the hash table?

Ans

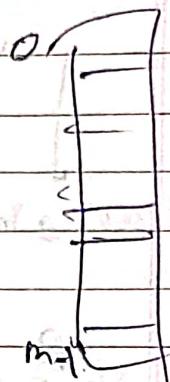
$$\text{Probability} = \left(\frac{1}{m}\right) \times \binom{n}{2} = \frac{n(n-1)}{2m}$$

collision can occur

at any one of  
the m places

↓

any 2 no.  
that we  
choose can



$H(x_1) = H(x_2)$  collide there

so if we take the size of the table in order of  $n^2$ , then

$$\text{Probability of collision} \leq \frac{1}{2}$$

but size of order  $n^2$  is too much. B/c in all the examples, we saw that the assumption for the amount of space that we can use is of order n.

→

A universal hash fn can be given in the following manner (no need to know as to why it is a universal hash fn)

$$H(x) = ((ax+b) \bmod p) \bmod m$$

$a, b, p$

$m \approx 2n$  (of order  $10^6$ )

$p \gg m \rightarrow$  of order  $10^9$  or greater

$$2 \leq 0 < a, b < p$$

We'll use all these things to create, what is known as the perfect hashing.

### PERFECT HASHING

A hash fn<sup>n</sup> is said to be a perfect hash fn<sup>n</sup> if there is no collisions at all.

Well in advance, if we know the no. which we want to insert into the table then we can create nicely hash fn which will help us to create the perfect hash fn.

We have a simple algo. which includes 2 level hashing.

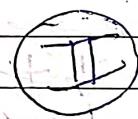


$n_i = \text{no. of elements which will hash to the } i\text{th location.}$

$n = \text{no. of elements to be hashed}$

$$m = 2n$$

$$\Rightarrow \sum n_i = n$$



The probability for choosing hash fn which will satisfy this property  $> \frac{1}{2}$

We first compute  $\sum n_i^2$ ; & if

Summation of space used in the secondary level

$$(\sum n_i^2) \leq 4n$$

then for the chosen hash fn

on avg. we may have to try with 2 hash fn

(Known as the primary hash fn.)

then the chosen hash fn<sup>n</sup> is correct, otherwise choose another hash fn<sup>n</sup> by choosing another P, a & b


There is another hash table at each location which will deal with  $n_i$  no. The size of this hash table is  $n_i^2$ .

That moment you tell me what value of  $n_i$  is  $n_i^2$ , I can choose a hash fn uniformly at random b/c you just use the following.

$$H(x) = ((ax+b) \bmod p) \bmod n_i^2$$

So instead of this use

$$\bmod n_i^2$$

So, we already saw that if  $n_i^2$

then  $\rightarrow$  Probability of collision  $\leq 1$

in 2nd level

hash table

there also  
on avg.

we may  
have to try  
2 hash fn.

$$\frac{1}{2} + \frac{1}{2} = 1$$

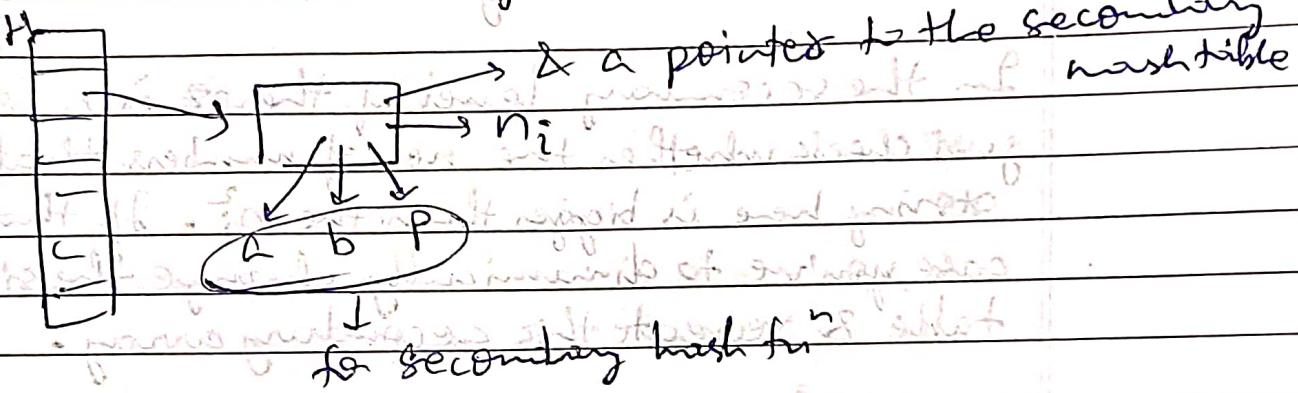
However, if there is still collision in the 2nd level, we choose another hash fn<sup>n</sup>.

1st level hash table is primary hash table  
2nd " " secondary hash table

i.e. the primary hash table will store the pointers to the secondary hash table. The second - any hash table will have <sup>different</sup> ~~another~~ hash fn.

In the secondary level, whatever we're going to write is not a linked list. Each of them is a hash table. Hash table will have a fixed length i.e.  $n^2$

So each index of primary hash table will have the following:-



For different secondary hash table at each index there <sup>may</sup> be different secondary hash fn

So, here the total space used will be

$$2n + 4n^2 = \leq 6n^2$$

∴ the complexity for finding hash fn is  $\Theta(n)$   
for both level

This is all what we do, when we're given the set of no. is already known to us. First we build this table & keep it. That is ~~when~~ there're no entries but we

know what are the hash fn's we're going to use. Now, if the entries are given, we can guarantee that there'll be no collision.

Q. But what to do if S is not given?

Sol: We start doing the same thing. There're 2 invariants needed to maintain. One is that the  $\sum n_i^2 \leq 4n$  & the 2nd thing worth mentioning is that at secondary level there is no collision.

In the secondary level if there's a collision, just check whether the no. of numbers that you're storing here is bigger than the  $n^2$ . If that is the case, you've to dynamically change the size of the table & recreate this secondary array.

The other thing that can go wrong is that  $\sum n_i^2 > 4n$ . In that case you've to change your entire hash fn from scratch. So, what you need to change is, ~~that~~ you've to change your primary hash fn. The moment primary hash fn changes, things will change in the secondary hash fn also. Whatever you've stored, you've to delete everything & then add it again. So if we don't know anything about S, this is the best that we can do.

So, in most of the time insert will take only order 1 but occasionally, when you've to recreate the hash table from scratch, it will take order n. So you have to resize the hash

Date: / /  
Page No.:

table & you also have to pick different hash table for different level. But that will not happen very frequently. that will take linear time  
Maybe once in a while it happens. So if you do n operations, only some of them will be order n. So overall it's expected that the complexity will be  $O(n)$ .