Data Structure

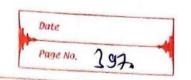
8. Hashing

Handwritten by pankaj kumar

Hashing

Page No. 396.

	Page No.	396.
	Introduction	(397)
0		(398)
2		.(338)
3	chaining	(394)
-	Open Addressing.	-
(5)	Harning in (++ 45ing STL	(401)
<u>6</u>	Two syn problem	(401)
(F)	4 sym	(402)
(8)	Longest Consecutive sequence in an array	(404)
(3)	Length of the longest sybarray with zero sum	(405)
0	Count number of subarrayer with given Kork	(406)
	Longest sybisting without any Repeating Character	(408)
	Congest 1900 101119 WILLIAM CITY KEPC STITY CHANGE	
		+-1



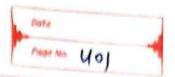
0	Hashing Totodyction in
	· Harring is a mergod or storing and tetrieving data from a
	databare efficiently.
εg	- For example we want to store employee record toyed using above
	numbers then we can use following data structure.
	(i) An array of phone numbers & records
71	(ii) A linked list of phone number & new rdg.
	(iii) A balanced binary search free with phone no. askeys.
	(i'y Adjrect access table.
	(i) Using an Array: & linked list
	searching of record: - O(n) or O(10gn) (cutost) (osty operation)
	Invert & delete: - O(1)
	(iii) balanced bingry search tree:
	moderate segret, insert delete: within O(logn)] = (moderate)
	(iv) direct Access table:
	emaking a big array and use phone numbers as index in the array
	· cue can do all operation in:- O(2)
	Souther tree when at the different control of
	· but this solution has many practical limitation, extra space required
	is huges, if phone no sigit is no then we need O(m * 100) space
	where mis size of pointer to the record.
Garage I	Another problem a programming language may not store 'n' digits.
£ 4	De Hadring a state in the state of the
	Hashing: Hashing is an improvement over Direct Access Table.
100 1	The idea is to use a hash function that convent a given phase
	number or any other key to a similar number and yser the
	small number as an index in a table called a hash table.
	19 Heat Fynetton: - A Function that converts a given by
	number to a small practicely integer value. The mappet integer
	value is used as an index in the hash table.

	· A good hash Function should have following properties:
	(i) It should be efficiently computable
	sin 11 should uniformly distribute the keys
	(Each table position be equally likely for each key).
	@ Had toller it moved conserved
	Howh table: An array that stores pointers to record consespond.
	ing to a given (phone number) An entry in hart table is NIL
	if no existing (phone nymber) has hash function value equal to the lindex for the entry.
	the Malex for the entry.
	in the first with the property of the property of
(2)	Collision Handling!
	. since a hash function gets us small number for a big tey.
	there is a possibility that two keys result in the same valve.
Arryl, elv	The situation where a newly inserted key maps to an already
	occupied slot in the hash table is called collision and myst
N. W. W.	be handled young some collision handling technique.
	collision handling techniqu:-
political	(1) Chaining: - The idea is to make each rell of the hard
	table point to a linked list of records that have the same
the state of	Function value chaining is simple but it required additional
The sale	memory outsides the table.
	(ii) Open Addressing: - in open addressing, all element are stored
and a	in the hash table itself. Full table entry contains either a
	record or a Nil, when segoiting for an element, are one by one
	examine the table slots until the derived element is tound no it
1	is clear that the element is not present in the table.
3	Chaining:
8 3	go to linked list section page no. (260-284)
- (4)	Q. Design a harbert outflow using byiltin hardtable
	Design of Hartmap without using any builting haryfalle

1) Open Addrewing:	
· like correcte chaining, Open addressing	is a method for
handling collisions in other magicialia	all planest and
in the hash table itself. so at any point	at the size of the
I IND CARDITION FINALL CONTROL FOR LA	
(Note: - we can increase table the b	Total number of keyr
(Mote: - we can increase table spe b	needed)
- Important operation:	
(i) insert(k): keep probing until an empt	h ol Lie C
an empty slot is found,	interest is found. Onee
(ii) Segrah(k): keep probing until the slot	in but k.
pand by by any the slot	to key doern't become
equal to k or an empty	Stot is reached.
(iii) Delete (k): Delete Operation is interesting	g. If we simply delete
a key, then the segges may	fail. So slots of the deleted
key goe marked specially	
Note: insert can insert an item in a d	eleted slot, but the
search doern4 stop cot a deleted slot	Ļ,
Open Addressing is done in the following	On Algue
	rig 4493 :-
① Linear Probing:-	1 the sale of the
IF slot haster) 1. S is full, then we t	14 (had(x)+1):1.5
if (hash (x)+1).1.5 is also full, then u	1e ty (hard (x) +2)-1,5
IF (h9sh (x)+2) 1.5 is also full, then a	ue to (bash(x)+3).1.5
The reader of the	-9nd 5000
(n's	nd an lard
	00
30117-1 (Mark)	50
700-1, 7 = O(19dex) /76-1, 7= 6	85
8 1.7 = 1 (colling an) : 1+1=2(idn) 5	32
92.1.7= 1 (cilliston), 1+1=2 (collisto):3 4	73
731.7= 3 (colling): 3+1=4	76
10/1.7= 3,:3+12v: 4+1=5	76

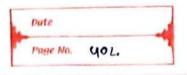
21.19

clustering: The main problem with linear probing is
clustering, many consecutive elements from groups and it
start taking time to find a free slot or to segoth an element
@ Quadratic Probing: - We look For (12 th) slot in (ith) iteration.
if (hash(x) 1.5) is full, then we try (hastx) + 1*1).1.5.
if (hash (x)+1*1) xs) is full, then ye try (hash (x)+2*2) 1/6
if ((hash(x)+2+2) 1.5) is full ; then we too (hash(x)+3+3).15.
a And Jo. on.
3 Double Hacking
3 Double Hashing: - we use another (hash 2(x)) function and
(or for 1 - harh 2(x) slot in ity italation.
· if slot (hash (x).1.5) is full, then we try (hash (x)+ 1+harras) is
" " (hash(x)+1* hash2(x)).1.5 " " " (hash(x)+2*hash2(x)).5
[(14) n(x) + 2 - 145 h 2(x)) 15
" [(h95h (x) +2*har2(x)) 1.5] "" " [(h95h(x) +3* h95h2(x)) 1.5].
(A) Conserving a of al
(*) Comparison of above three
· linear probing = best cache performance but suffer from clustering
and early to compute, boost reacher morting more
· Quidratic possing = lies blug the true in terms of cause & clustering.
Double performance but no chotenno.
Double harbing require more computation fine.
Performance of Open Addresting:
m= Humber of slots in hash table in = number of fey to be inserted
Logd Factor X - Ol Free
[Load factor of = n/m (<1)
Expected time to search insert (selecte < 1/(1-4)
So, Segocy, Net/ Setell take [1/(1-4)] time

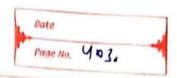


1	Hashing in Ctt Using STL
	. Hashing can be implemented using drop
	present in STL as per the requirement
	CD Set
	(ii) Unordered set
	(iii) map
	(iv) Unordered may
	Mote: - [all the topics covered in 17+572 page no. 93-100]
(6) Tup Sym
	· Given an array of integers nums and ntarget, return the
	indices of the two number such that they add up to the target.
	EX nums = [2/7/11/15], target = 9, off [0,1]
	Ελ 11411 (1/2) 14/9 -3, 0/ρ (0,1)
	Approach 1: 4ring two loop: 0(n2)
	Approach 2:- Optimized (Hary-table)
	· go through all the element take current element and check its
	other pair are present or not in high-table if not present than
	Project it into hash-table and go through next elements
	Code:
	vector (int) two Sym (vector (int) & nymr, int taget)
	vector-lints ans;
l.	unordered_map < int, int> mpp;
	For (int izo; ic nyms. size u; i+t)
	if (mpp. find (target-numscia) ! = mpp. endu)~
	ans. push - back (mpp farget - numr tist);
	ans. puth-back (i);
	s return ans; 111
)——	s mpp Trymstiz]=1: //TC:06)
	return 905: [/ S(:010)

(0.4)

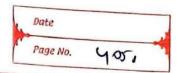


(1)	4 sym	
	Given an array nyms of n integers and an	
	are there elements a by G and In nymr syc	4 that.
	9+b+(+d= target? Find all unique gyadryp)	ets in the
	array which gives the um of target.	
	[r:- nums=[2,0,-1,0,-2,2], and targ	et=0.
	TOPE T-	
	<u>off</u> [-1,0,0,1],	
	[-21-111,2],	
	J [-2101012]	
	Approach 1:- Sort -> 3 pto + Bingry segrit -> 5	Cal Co
1.30		et (60 401948)
kangala di	51112233444 + 1918el=9	Market Market
	11 12 23 44 9	
	9-3=6)>85 to Find 6	1000
	if not found increment k,	iso
	11122 33 44 4	j= j+1 == k=j+1
Mark 1 and		- 1 - 11T
	9-4=5 -> Bis to Find 5	is i
	and so. on Find all quadruplet of push into	haphsetr
	T(: 0(N3 109N) + 0 (N109M) / 51:0	(1)
	Approach 2 optimized	
	O Sort Course course con the contract of	
174	2[112233444] fa	yet=9
. 615	à d'est	7
- (9	(9-11+3)]= (7 in (left to right)	
	[3 £1 2233444]	
	10FT +184+	
	: 1+4<7 (so we have to inv	
	1/10 de 1965e value 1.	e increase left)



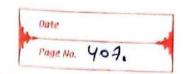
is tell left left left 4 4 4 4 584+
2+4<7 (increase left)
How sty = 2 to due and
(1)11811) Put it linto hash set ~ (1,1,3,4)
hashed repeat their if (! (eft >= Hight) of if (1eft < right)
Now increase 1
[1 1 2 2 3 3 4 4 4 7] j=1+1
Mour repeat step (2) Find (target-(1+5)) in [left, right]
2/22 / 2/2/
SC: 0(1) Note: don't need to use Hastref
- sole 'r
Trector (rector (int)) Foursym (rector (int) & humr, int target)
vector (vector (int>) total;
int n=nums. size();
if (n<4) retym total;
sort (nums. begin (), nums. end();
for (int 1203 icn-3; 1++) of
if (i>0 2f nymsti]== nymsti-1]) continue;
if (4*nums ci') > +arget) brokeak;
For (int 5= i+ 1; i< n-2; i++)
if (5) i+1 23 nyms (1) == nyms [5-1]) (on thye?
if (nums [i] + 3* nums [i] > target) break;
int left= i+1, night=n-1;
while (left < right)
int symenums tieft Jthums Dight Jthums Cistowns
if (sum < target) left + ti
else if (sum > target) right;
total push back (vector (int) & nymscis, nymscis
total. push-back (vector (int) of numscis, numscis)

Date lefteright).



	Page No. 405,
	longeststreak = max (Tongeststream, cyrrstream);
	e return longest Streak; / 1000
	3
(9)	length of longert sybgingy with zero sym:
(£x	2N=6, arco= 29,-3,3,-1,6,-55, 0/P= 5
(%)	sybamy sym with zero are.
	8-313518-1,61-56, 8-3,3,-4,6,-55.
	longert is = L-5,13,-4,6,-5}, Length = 5
E	N=5, arc 3= 21131-5,6,-25, 0/p20
	Here is no subgray that sums to zero.
_	Approach 1:
	Had longest subgroup with memor by item was
	100 por 100p:- for (100 to ich)
	for (j=1 to j'(1)
	$T(:0(N^2))$ $Symt = arcij$
	Operand of (sum=20) (alculate length (j-i+1)
	Approach 2:
	that a close the prefix sym of every element, and if we observe
ę.	that 2 elements have some prefix sym, we can conclude that the
	Ex: arr = [9, -3, 3, -1, 6, -5] rice subgray (i,k) = subgray
X 1	prehirsum = [9,6,9,8,14,9] = then, (Ex); sym is zero
	Symis zero (3,5) f (2,5) have sym 15 zero
- 100	and here length of (1/5) is maxim so. this will be answer,
	for this we can use Hart man of there index with now in
	T TOUR LEGICAL DISTRICT THE PROPERTY OF THE PARTY OF THE
	check if it is of any prefix sym from hassmap, if we find zero sym then it will also be the salisfies consum
	The sen syn then it will also be the salisfied contin

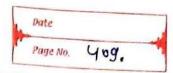
int maxlen (int ACT, into) unordered-map (int) int> mpp; int maxi = 0; int sym= 0; for (intiosicn si+t) of SYM +=ACis; if (sym==0)~ Maxi = 1+1; cired if (mpp. And (sym) ! = mpp. end())} Maxi = max (maxi, i- mpp Gym); elsed mpp crum = i; Fetyon maxis Count the number of subarrays with given xork Given an array of interger A, and an integers. And the total number of subgarays having bitwise xor of all element equal to B EN:- A = [4,2,2,6,4] 1826 explantion sysappray having to or of all element 6 gae [4,2] 5[4,2-12,6,4], [2,2,6],[6] Approach 2 > TRRUTE FORCE II O(N2) generate all possible sybarray for each sybarray get respective XOR and check if it is equal to B then increment grower count to 1. in the end we will get count of subbroan with xor egyal to



	Approach 2: (prefix xor and map):-
	s prefix xox megar xox from thex o to that olement.
	mate a garay of prefix too for all olar
	prefix=x00 Ei] = X0R(aco], 9 [1], 9 [2], 9(i))
	Note: Obsaration:
	P= x0x (9 507, 0527, 9507, 9597, 959+17, 9 TP])
	Q= x0x [960],961], 969])
	: p1 8 = x08 (9 [9+1], - 9 [p]) = 17.
	so, we understand that from xor grossy when we xor two
	element out different indices we get the xor of the element
	(in the original array) blu those two indices.
	Mote:
<u>a</u>	let's say we did XOR (PIB) and we get Q (Bir given integer)
	this means that the subgroup blu q and p having now
	=B. exi- pnB=Q
	=> P^B^P= Q^P
	=) B=QAP
Steps	This means we will iterate over the good god to Revery element
	and check if it is B then could be incremented by 1.
Ste	P2: - we will also store every-prefix_xor with frequency in
	map.
<i>A</i>	13:- Yor the current prefix with B (i.e PAB, where, P= 965) 965-765
N. S.	if PMB=Q is present in map. it means from 8+1 to P
Land.	there is suggence with NOR B. now increment the count by
1/5	Foregrency of 8.
	TC:-ON)
	SC:- O(N)

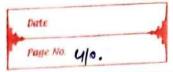
FLO

	Co4e:-
	# include < bits/stdc++.4>
4000	using namespace std;
	class Solutions
	public:
(11)	int solve (vector (int > PA, int B)
	unordered-map cintint's visited;
2.761	int cpx =0: (code)
5	long long coo;
	For (int i=0; i < A size(); i++)d
	$CPX^{\wedge} = ACIT$:
	IF (CPX==B) C++;
0.641 9	10+ 4= CPX B:
4-9-5	if (visited. Find (4)!= visited. end 1)?
	suisol out & C=C+ Visited[4];
	S Visited EcpxJ++;
and al	retum c;
	Programme and the state of the
	Vector < in+> AQ 4,2,2,6,4.5;
	Solution 065;
	int total Count = D.b.T. Colve (ALC)
	S Contretotal Countre en 20;
make to (
King Pal	Dength of longest substring without any Repeating character.
6.	Ep: I/p: 5=11ebcabcbb" /ojp: 3
	7/0: 5=11 bbbbb11 / op: 1
A. O. Y	Approach 1: Brute Force:
6 A	generate all the substring by taking For loop one by
7 23	one check for each a eurn element if the element is already
1	violted in the current substains then find the length and
	tetum break from the nested roop
	· to find if element it present or not in current substring one
	can use has 4 set



11	
	int solve (string str) of
	if (stor. size t) ==0) return a.
	INT MARGIE = I M-MTN!
	For (int 100; 1< stor lengths); (++) of
	unordered - set all lead
	100 (11/15/15/15/10/10/10/10/10/10/10/10/10/10/10/10/10/
	if (set. Find (stable)) = set. end())
	11 0x 1 min
	set. invertato Tis);
	telurn maxanr; Tr: O(N2)
	\$ \\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\\
	=======================================
	Approach 2: Ophimised solution 1 1:
and too	we will have two pointers left and right pointer left is
	ysed for maintaining the starting point of substring while
	Bight will maintain the end point of the sybstoing. Fight
	pointer will move forward and check for the duplicate ownnerse
	of the current element if found then left pointer will be shifted
	ahead on as to delete the duplicate elements.
	Eni- abragbedba abedba
	abcaahedha aheadha
	The contract of the contract o
	abcerabedba - abcqquedba
	2 (v)
	THE CLOUGHT OF THACH
	int solve (stoing sto) of
	if (shisizer) retum 0;
	int maxau = Int-wing
	unordered_set
	int leo;

* PT5



	For (int r=0; re str. length U; r++)~
	if (set find (sto cr1)!= set.endu) of
	while (ecr pg set. Find (str [r]) != set. end()
	Set. ergre (sto Cas):
	£ 1++;
-	set insert (stores);
100	5 maxans = max (maxans, 8-0+1);
4	1 (11/4/15/8-D+1);
	E telan wax an;
	Tr: 0(2*N) (sometimes left of right both have to travel
	(1: O(N) (Hashard of rice M) Complete array)
	Approach 3: Optimized Solution 2:
16511	eye will make a map that will take ase of counting the
1 100	elements and maintaining the frequency of each and every
211.40	element as a unity by fating the latest index of every element.
ME ST	Code:
(17 K)	int length of Longest substring (string s)~
	vector xint > mpp (256, -1);
. A tour	int left = 0, right = 0;
	1nt n= s.sive();
Low t	int len=0;
	while (right < n) ~
4 1 5	off (mpp [scrish]] 1=-1)
177	left=max(mpptscnight]+1, left);
	mpp [stright] = right;
	leh = max (len risht-left+2):
	6 818h+ 13 (mos) (33)
	May and an and an
	return len; Trio(N)
	\$ St:0(M)