

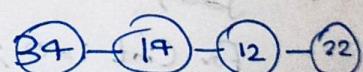
HASHING

The Problem

- There is a large Phone Company, they want to provide caller ID capability.
- given a phone number, return caller's name.
 - phone no. range from 0 to $10^8 - 1 = r$
 - There are n phone numbers, $n \ll r$
 - Want efficiency.

Soln 1: Unordered Sequence

- Searching & removing : $O(n)$
- inserting : $O(1)$
- Used when:
 - Frequent insertion
 - Rare searches & removal

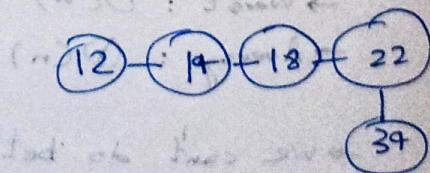


Soln 2: Ordered Sequence

- Searching : $O(\log n)$ [Binary]
- Insertion & Removing : $O(n)$

Application

- Frequent Searches
- Rare insertion & removal



Soln 3 : Direct Addressing

- An array indexed by key
- ~~the~~ addressing take $O(1)$ time., but huge amount of space required. $[O(r), r = 10^8]$

(null)	(null)	(null)	Ankur	(null)	(null)
0000-0000	0000-0000	0000-0000	96 35-8904	0000	0000

Soln 4 : Hash Table

→ Addressing Time : $O(1)$

→ Size Requirement : $O(n+m)$

↗
 no. of phone numbers
 ↗ table size.

→ Like an array, but come up with a function to map the large range into one which we can manage

Eg Function : Take original key, modulo the (relatively small) size of array, and use that as index

Eg : For Ankur, hashed into 5 slots

$$\text{index} = 96358904 \% 5 = 4$$

(null)	(null)	(null)	(null)	Ankur
0	1	2	3	4

Example

→ Let keys be entry no. of student in Batch

CS 302

Eg. 2018 CS 10 110

→ There are 100 students in class, so we can't create a hash table of size 100

→ Let Hash Function output last \rightarrow^2 digit

Eg: 2018 CS 10 110 \rightarrow 10

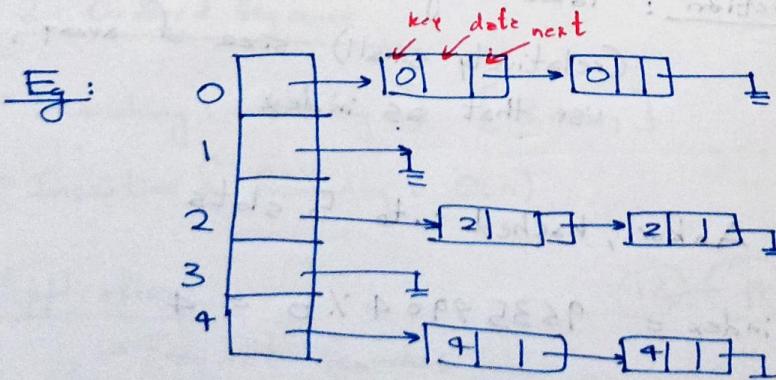
but also 2018 CS 50310 \rightarrow 10

Collision.

Collision Resolution

→ To Solve Collision, we can use chaining

→ Set up an array of links (a table), indexed by keys, to list items with same key



→ Most Time Efficient Collision Resolution Scheme.

→ To find / insert / delete an element

- $\xrightarrow{h(k)}$
- Using Hashing Function, Find at key $h(k)$
 - Using the key, you have a linked list,
do Linked List Find / Insert / Delete.

Analysis of Hashing

- An element with key k is stored in slot $h(k)$ [instead of slot k without hashing]
- The hash function h maps the universe U of keys into slots of hash table $T[0, \dots, m-1]$

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

- Let time to compute $h(k)$ be $O(1)$

- A good hash function is the one which distributes key evenly amongst the slots

- An ideal hash function would pick a slot, uniformly at random and hash the key to it. But this can't be used as we can't find the slot to use while searching for that element.

Simple Uniform Hash Function [SUHA]

- Hash Function is assumed SUHA when
 - i) function will evenly distribute item to slot
 - ii) each item to be hashed has equal probability for each slot. [regardless of situation of table]

→ Given a hash table T with m slots holding n elements

$$\text{Load Factor} = \alpha = \frac{n}{m}$$

To explain

Unsuccessful search

→ Element not in linked list

→ Simple uniform hashing yields an average list of length $\alpha [\frac{n}{m}]$

→ Expected no. of elements to be examined =

∴ Search time : $O(1 + \alpha)$

hash computing searching

Successful search

→ Assume that a new element is inserted at the end of Linked List

→ upon insertion of i^{th} element, expected length of the list is $\frac{i-1}{m}$

before insertion of i^{th} node

→ in case of a successful search, the expected number of elements examined is 1 more than the number of elements examined when the sought-for element was inserted.

(previous nodes + 1 required node)

Now to find

- For i^{th} element, no. of elements examined = $1 + \frac{i-1}{m}$
- For average ~~time~~ elements for whole table, we do

$$\begin{aligned} \frac{\sum_{i=1}^n \left(1 + \frac{i-1}{m}\right)}{n} &= 1 + \frac{1}{nm} \sum_{i=1}^n (i-1) \\ &= 1 + \frac{1}{nm} \times \frac{(n-1)n}{2} \\ &= 1 + \frac{n}{2m} - \frac{1}{2m} \\ &= 1 + \frac{\alpha}{2} - \frac{1}{2m} \end{aligned}$$

Time Complexity = $\Theta\left[1 + \frac{\alpha}{2} - \frac{1}{2m}\right]$

process = $\Theta(1 + \alpha)$ [generally $m \rightarrow \infty$ $m = \text{constant}$]

- IDEA: A slot of the table prints a set.
- Assuming the number of hash table slots is proportional to the no. of elements in the table.

$$n = O(m)$$

$$\alpha = \frac{n}{m} = \frac{O(m)}{m} = O(1)$$

∴ Searching takes constant time on average.

→ Insertion takes $O(1)$ time

→ deletion takes $O(1)$ time, if linked list are doubly linked.

Good Hash Function

- The function which can be computed quickly
- It should distribute the key uniformly over the hash table
- Very Rare
- Eg: Birthday Paradox
 - Even in a group of 35 or more students in the class there is a very high chance that 2 of the students would have same birthday.

How to deal with non-integer key

- Find some way to convert key to integer

Eg: $96\ 35 - 8904 \rightarrow 96358904$

→ for a string, add up it's ~~ASCII~~ ASCII values.

Hash Functions

- The mapping of keys to indices of a hash table is called a hash function
- A hash function is usually the composition of two maps, a hash code map and a compression map
 - An essential requirement of the hash function is to map equal keys to equal indices
- A good hash function minimises the probability of collision.

Hash code Maps

keys \rightarrow integer

Compression map

integer $\rightarrow [0, 1, \dots, m-1]$

Popular Hash Code Map

Integer Cast

→ for numeric types with 32 bits or less,
we can reinterpret the bits of the number
as int

Component Sum

→ For numeric types with more than 32 bits
Eg long & double, we can add 32 bit component
→ Bad for strings.

Polynomial accumulation

→ For strings of a natural language, combine the character values (ASCII or Unicode) $= z_0, z_1, z_2, \dots, z_n$
by viewing them as the coefficient of a polynomial

$$S = z_0 + z_1 \cdot n + \dots + z_{n-1} \cdot n^{n-1}$$

→ Horner's Rule :

$$S = z_0 + z_1 \cdot n + \dots + z_{n-1} \cdot n^{n-1}$$

[can be done by recursion]

$n = 33, 37, 39, 41$ is observed with at most 6 collision on a vocabulary of 50,000 words

Compression Maps

(full class notes)

$$\text{I) } h(k) = k \bmod m = k \% m$$

$k \rightarrow \text{key}$
 $m \rightarrow \text{size of hash table}$

$m \rightarrow b^e \text{ form} \rightarrow$ bad. slow. death. robust
 if m is power of 2,
 $h(k)$ gives LSB of k
 all key, ending with same
 ending go to same place.

prime → good
 ensures uniform distribution
 prime should not be too
 close to powers of 2

Eg: $n = 2000$ character string.

→ we can handle 3 nodes in LL

∴ we can do $m \approx \frac{2000}{3}$ but prime

$$= 701 \text{ (not near power of 2)}$$

$$\text{II) } h(k) = \lfloor m (kA \bmod 1) \rfloor$$

$k \rightarrow \text{key}$

$m \rightarrow \text{size of table}$

$A \rightarrow \text{constant, } A \in (0, 1)$

Steps involved:

- map $0, 1, \dots, k_{\max}$ to $0, A, \dots, Ak_{\max}$
- take it's fractional part $\bmod 1$
- map fractional part to $0, 1, \dots, m-1$

→ value of m is not critical, generally use $m = 2^P$

→ value of A depends on characteristic of data.

$$\text{Fibonacci Hashing: } A = \frac{\sqrt{5} - 1}{2} \quad [\text{by Knuth}]$$

conjugate of golden ratio.

$$\text{III) MAD [Multiply, Add, Divide]}$$

$$h(k) = \lfloor ak + b \rfloor \bmod m$$

$k \rightarrow \text{key}$

$m \rightarrow \text{size of hash map}$

$a, b \rightarrow \text{const.}$

→ eliminates pattern, given a not multiple of m

→ same formula used in linear congruential (pseudo)
random number generators

Universal Hashing

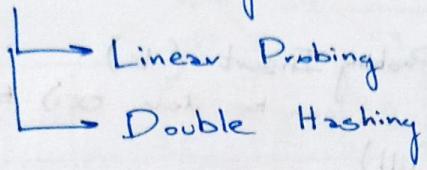
- For any choice of hash function, there exists a bad set of identifiers
- A malicious adversary could choose keys to be hashed such that all go to same slot.
- Average Time Retrieval becomes $\Theta(n)$
- Solution:
 - A random hash function
 - Choose hash function independently of keys!
 - create a set of hash function H , from which h can be randomly selected
 ↓
 h should be same in a runtime
- A collection H of hash function is universal if for any randomly chosen f from H ,
(and two key k & l),

$$\text{Probability } [f(k) = f(l)] \leq \frac{1}{m}$$

Collision Handling Techniques

→ Chaining [LL method]

→ Open Addressing



Open Addressing

→ All elements stored in hash table [can fill up]

i.e. $n \leq m$

→ Each table entry contains either an element or null

→ When searching for an element, systematically probe table slots.

→ Modify hash function to take the probe number i as the second parameter.

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

→ Hash function h determines the sequence of slots examined for a given key.

→ Probe Sequence for a given key k given by

$$\{h(k, 0), h(k, 1), \dots, h(k, m-1)\}$$

⇒ permutation of $0, 1, \dots, m-1$

Linear Probing

→ If current location is full used, try next location

Algorithm Linear Probing Insert (k)

can be done $O(1)$ by counter

if (table is full)

error

probe = $h(k)$

while (table[probe] occupied)

probe = $(\text{probe} + 1) \bmod n$

table[probe] = k

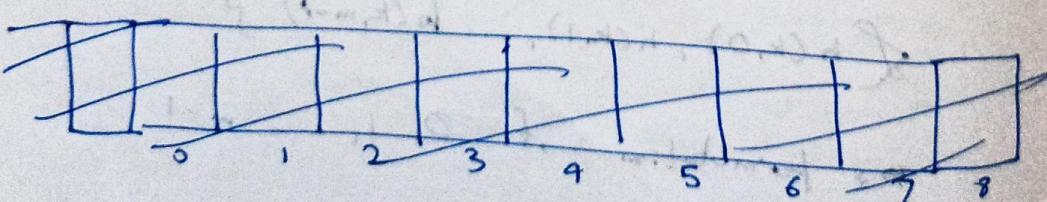
}

→ Uses less memory than chaining as no links or LL

→ slower than chaining since one might have to walk along the table for a long time.

Eg: $h(k) = k \bmod 13$

keys $\rightarrow 18, 41, 22, 44, 59, 32$



Eg: $h(k) = k \bmod 13$ [table size = 13]
 keys : 13, 41, 22, 44, 59, 32, 31, 73

				31					73		
				44	32						
				18		59					
					18	44	59	32	22		
									22		
									31		
										73	
0	1	2	3	4	5	6	7	8	9	10	11
											12

Look up in Linear Probing [Searching]

→ To search key k , we go to $h(k)$ and continue looking at successive location till we find k or encounter empty location.

→ Successful Search:

→ Eg. 31

$$h(k) = 31 \bmod 13 = 5$$

check loc 5 → 18 X
 6 → 44 X
 7 → 59 X
 8 → 32 X
 9 → 22 X
 10 → 31 ✓

∴ 31 is present at loc 10

→ Unsuccessful search.

→ Eg. 33

$$h(k) = 33 \bmod 13 = 7$$

check loc 7, 8, 9, 10, 11 don't find
 loc 12 → empty → 33 not present.

Deletion in Linear Probing

- Search number and delete it
- suppose in previous example we deleted 32
we get

		41		18	44	59		22	31	73	
0	1	2	3	4	5	6	7	8	9	10	11

→ Now if we search - fro 31

$$h(31) = 31 \bmod 13 = 5$$

check 5, 6, 7 not found

8 empty → Not present
(false)

→ Solution:

→ instead of placing bss 8 empty [null], place
→ tombstone [marker]

→ lookup should ignore tombstone and proceed
normally

→ If insert comes at tombstone, remove tombstone
and insert there.

→ Too many tombstones can degrade lookups performance
Rehash if too many tombstones.

Double Hashing

→ Uses two hash function h_1, h_2

$h_1(k)$ → position where we first check for key k

$h_2(k)$ → determines offset we use when searching for empty position

Algorithm Double Hashing Insert (v)

{ if (table is full)

 error

 probe = $h_1(k)$

 offset = $h_2(k)$

 while [table[probe] occupied]

 probe = (probe + offset) mod m

 table[probe] = k

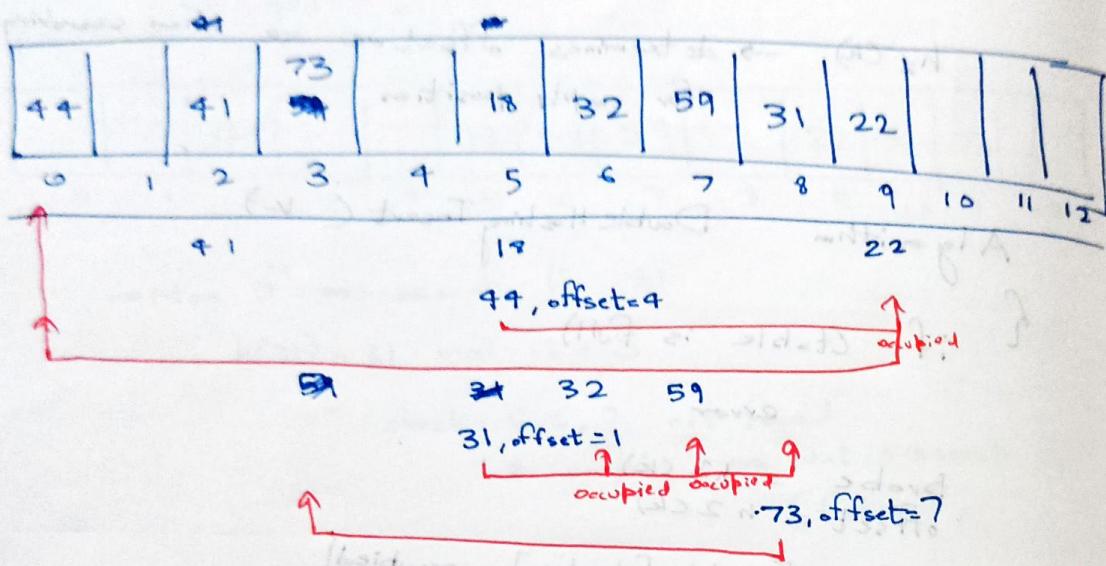
}

→ m should be prime ~~to~~ for covering every position in table with any offset.

→ Same disadvantages as Linear Probing.

→ Distributes key more uniformly than Linear Probing.

Eg: $h_1(k) = k \bmod 13$ [size 13] slots
 $h_2(k) = 8 - (k \bmod 8)$
 keys: 18, 91, 22, 49, 59, 32, 31, 73



Analysis

→ we assume that every probe looks at a random location in the table.

→ $1 - \alpha$ is the fraction of empty table.

MERGE SORT

Divide and Conquer

- Divide the problem into a no. of sub problems
- Similar sub problem of smaller size.
- Conquer the Sub Problem
 - solve the subproblem recursively
 - sub problem size small enough: solve problem in straight forward manner
- Combine the solutions of the sub problem
- Obtain the solution for the original problem

Merge Sort Approach

I) Divide

- divide the n -element sequence to be sorted into two subsequences of $n/2$ element each

II) Conquer

- Sort subsequence recursively using merge sort
- when the size of subsequence is 1 there is nothing more to do

III) Combine

- Merge the two sorted subsequence.

E. jorithm Merge Sort (A, l, r)

```
1 if ( $p < r$ ) // Check base case  
then    $q \leftarrow \lfloor (l+r)/2 \rfloor$  // Divide  
        Merge Sort ( $A, l, q$ ) // Conquer  
        Merge Sort ( $A, q+1, r$ ) // Conquer  
        Merge ( $A, l, q, r$ ) // Combine
```

{

Initial call: Merge Sort (A, l, n)Eg: n power of 2

1	2	3	4	5	6	7	8
5	2	4	7	1	3	2	6

$q = 4$

1	2	3	4
5	2	4	7

$l=2$

5	6	7	8
1	3	2	6

$q=6$

1	2
5	2

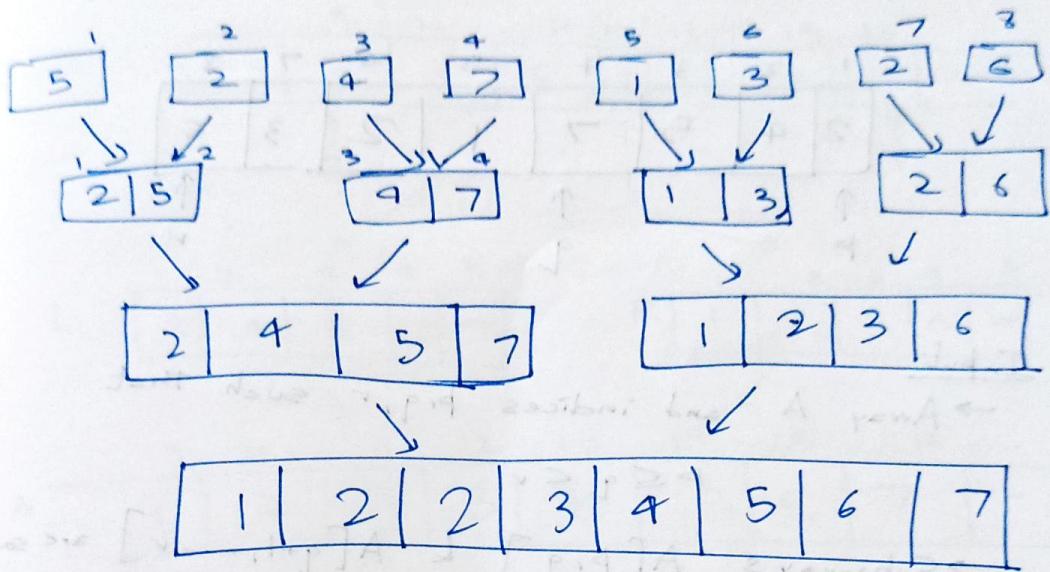
3	4
4	7

5	6
1	3

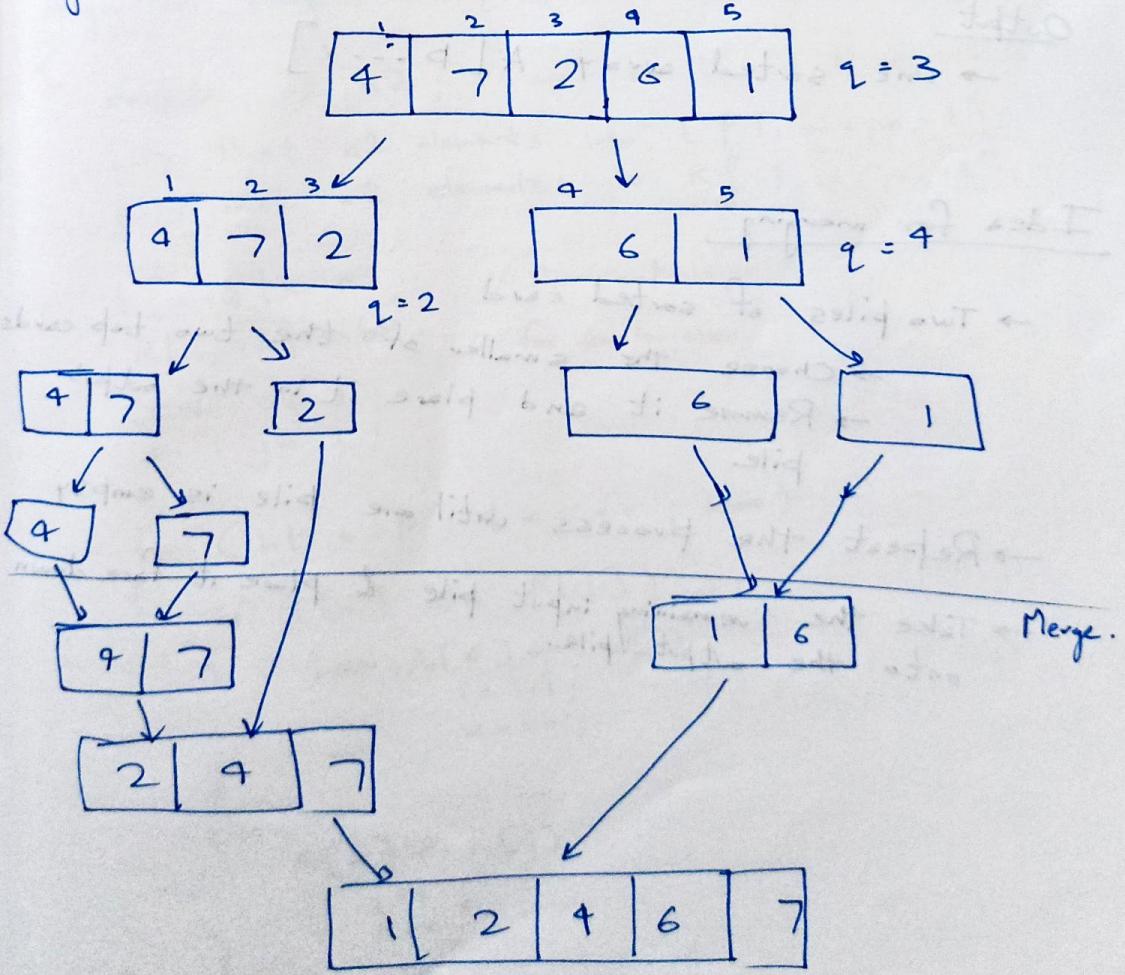
7	8
2	6

7	8
2	6

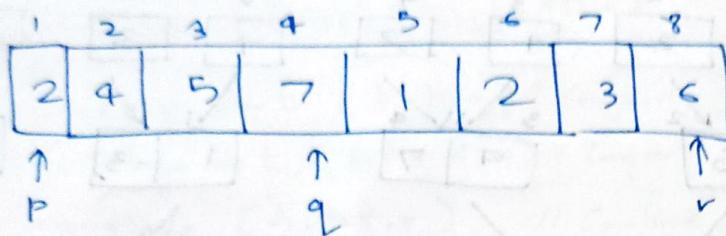
conquer & merge.



Eg. not power of 2



Merging (very similar to merging we did in polynomial addition, sparse matrix addn...)



Input

→ Array A and indices p, q, r such that

$$p \leq q < r$$

→ Subarrays $A[p, q]$ & $A[q+1, \dots, r]$ are sorted

Output

→ One sorted array $A[p \dots r]$

Idea for merging

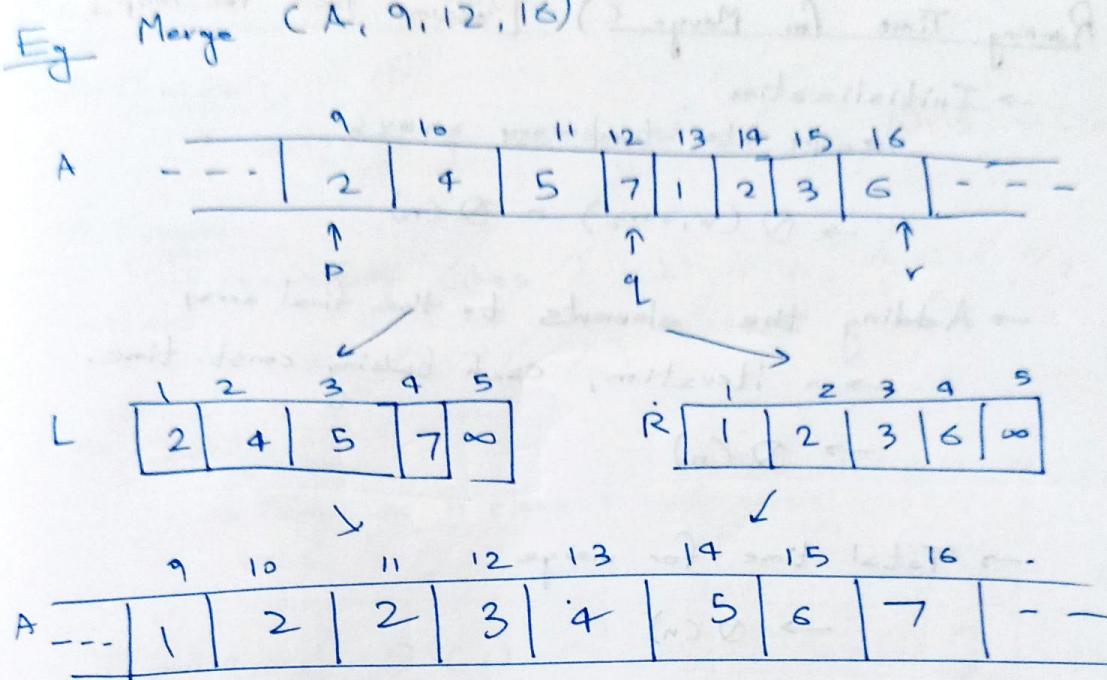
→ Two piles of sorted card

→ Choose the smaller of the two top cards

→ Remove it and place it in the output pile

→ Repeat the process until one pile is empty

→ Take the remaining input pile & place it face down onto the output pile.



Algorithm Merge(A, p, q, r)

```
{
    compute n1, n2; // size of L and R array
    copy first n1 elements into L[1, -- n1+1]
    copy next n2 elements into R[1, -- n2+1]
```

$L[n_1+1] \leftarrow \infty$; // for comparison

$R[n_2+1] \leftarrow \infty$; // for comparison.

$\exists i : i \leftarrow 1$
 $\exists j : j \leftarrow 1$

for (int k = p; $k \leq r$; $k++$)

{ if ($L[i] \leq R[j]$)

{ $i++$ $A[k] \leftarrow L[i]$
 $i = i + 1$
 $k = k + 1$

}

else {

$A[k] \leftarrow R[j]$
 $j++$
 $k++$

}

} add all ~~values~~ remaining values as it is

Running Time for Merge() [assume fast for loop]

→ Initialization

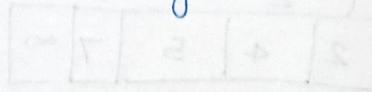
→ copy into temporary arrays

$$\rightarrow \Theta(n_1 + n_2) = \Theta(n)$$

→ Adding the elements to the final array

→ n iteration, each taking const. time.

$$\rightarrow \Theta(n)$$



→ Total time for merge.

$$\rightarrow \Theta(n)$$

Analysing Divide and Conquer Algorithms

→ The recurrence is based on the three steps.

→ $T(n) \Rightarrow$ running time on a problem of size n

→ Divide ~~the~~ the problem to a subproblems, each of size n/b : take $D(n)$ time.

→ Conquer (solve) the subproblem, a $T(n/b)$ time.

→ Combine the solⁿ $C(n)$

$$\therefore T(n) = \begin{cases} \Theta(1) & n \leq c \\ aT(n/b) + D(n) + C(n) & \text{otherwise} \end{cases}$$

Merge Sort Running Time.

→ Divide:

→ computes q_L as avg of q_L & r : $\mathcal{O}(1) = \mathcal{O}(n)$

→ Conquer:

→ recursively solves 2 subproblems, of size $n/2$

$$\Rightarrow 2T(n/2)$$

→ Combine:

→ Merge on n element subarray : $\mathcal{O}(n) = \mathcal{O}(n)$

$$\therefore T(n) = \begin{cases} \mathcal{O}(1) & n=1 \\ 2T(n/2) + \mathcal{O}(n) & n>1 \end{cases}$$

~~$\therefore T(1) = \mathcal{O}(1)$~~

~~$T(2) = 2T(1) + \mathcal{O}(2)$~~

⋮

~~$T(n-1)$~~

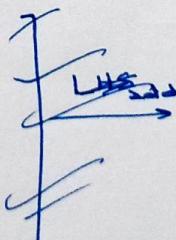
~~$T(n) = 2T(n/2) + \mathcal{O}(n)$~~

$$\therefore T(n) = 2T(n/2) + n$$

$$T(n/2) = 2T(n/4) + \frac{n}{2}$$

⋮

$$T(1) = 1$$



no. of levels : $\frac{n}{2^k} = 1$

$$n = 2^k$$

$$k = \log_2 n.$$

$$\therefore T(n) = 2T\left(\frac{n}{2}\right) + n \quad n \mid$$

$$T\left(\frac{n}{2}\right) = 2T\left(\frac{n}{4}\right) + \frac{n}{2} \quad \times 2$$

$$T\left(\frac{n}{4}\right) = 2T\left(\frac{n}{8}\right) + \frac{n}{4} \quad \times 2^4$$

:

($\frac{n}{2^k}$) T

$$T(1) = 1$$

$$\times 2^{\log_2 n} = n$$

add.

$$T(n) = n + n + n + \dots - \underbrace{\dots}_{(k \text{ time})}$$

$$= nk$$

$$= n \log_2 n$$

$$= n \log n$$