



G

O

D

O

O

O

O

D

Disjoint set Data Structure

/ Union-Find Data Structures

— as an abstract data type (ADT)

— implementation so that it's efficient for my algo

ADT

Data type

integers -

uses defined
data type

** In some cases, actual data type
is not so important

↙
it helps in my
algo

Example :-

For your algo, you need a

FIFO queue of items.

Q need certain operations

→ insert items in queue

→ remove items from queue (FIFO)

→ Is the queue empty?

FIFO
Queue
ADT that supports
these operations

It is convenient that we agree on the operations

without specifying the data type of the items.

& design my algo

Define a data structure \rightarrow Implementation view



(Not the only view)

operational view ✓

implementation view ✓

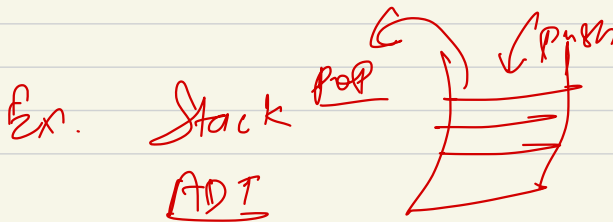
FIFO queue: insert, remove...

┌ - arrays
└ - linked list

User of the data structure doesn't need to know about implementation view.

Operational view \rightarrow A general ABSTRACT DATATYPE

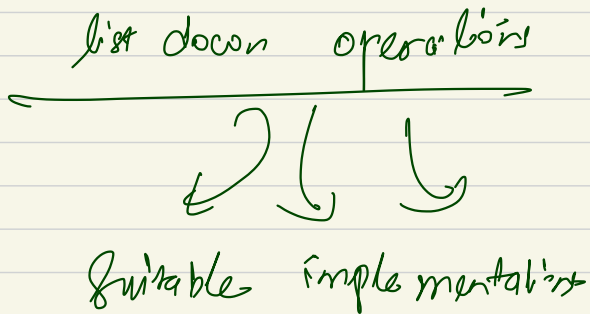
more general



Operations

┌ implementation
└ - arrays
└ - linked list

While designing an algo, if you realize that your needs correspond to a part^o ADT you immediately use that.



Modular process

Algo

uses the ADT (without worrying about implementation details)

* We cannot always ignore these.

Disjoint set / Union find as an ADT.

↓
nice example of not so straight-forward data structures that helps in improving efficiency of Algos

Settings ⇒ There are n elements x_1, \dots, x_n

divided into disjoint sets

Initially, each element is a group in itself

Operations ⇒

Disjoint set
data structure

MakeSet(x) ⇒ Make a singleton set with x id = "x"

FindSet(x) ⇒ Given x, find the id of the set it belongs to
or Find(x)

Union(x, y) ⇒ Create a single set of two sets containing x & y. New set can have id of any of two

Example of an algo:- (Min Spanning Tree Algo)

Kruskal's Algo

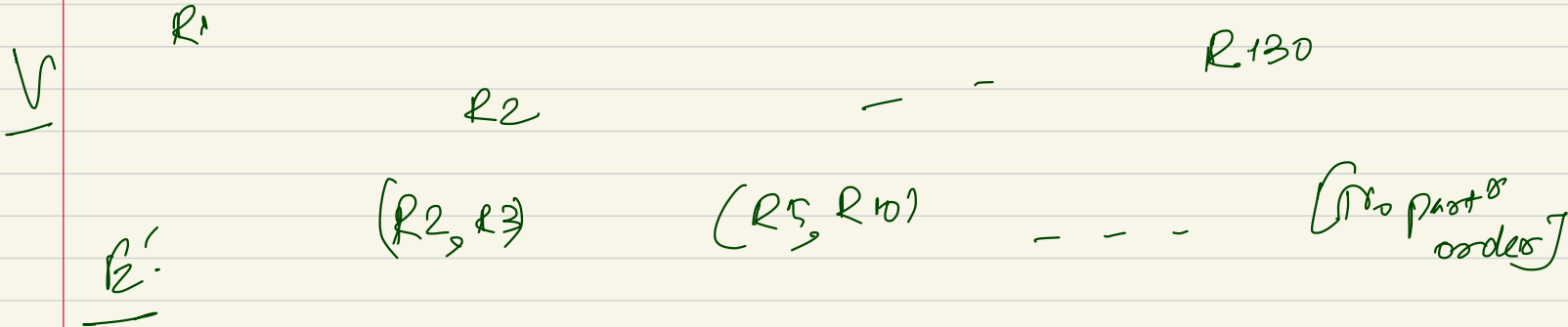
Setting

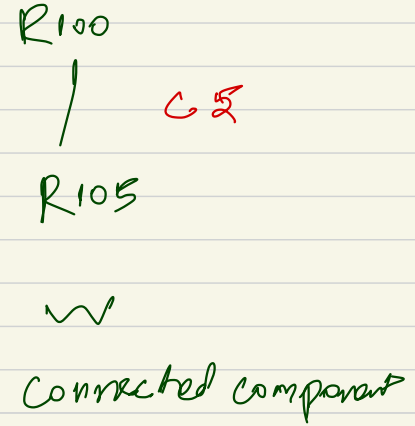
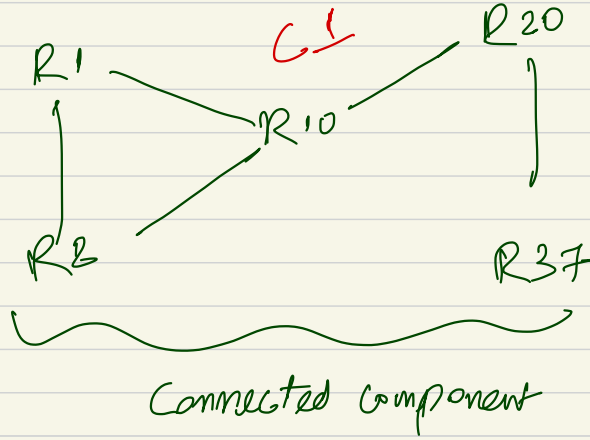
Think of a social network GN , (Facebook/Twitter)

defined by a set of people V

& a set of edges that define

relationship betⁿ pair of people E (in no part^r order)





Write an algo that finds out all the connected components,
and for each person $x \in V$, output the connected component
 x belongs to

R2	C1
R100	C5
⋮	⋮

* You can't make any
assn about order of E
(o/w, traversal could
be used)
- You get to scan E only once

Algo (that uses my ADT)

for each $x \in V$ do
 makeSet(x)

for each $(x, y) \in E$ do
 if find(x) \neq find(y)
 union(x, y)

for each $x \in V$
 output - Person " x " belongs to component
 find(x) .

m edges n people

$2m$ - find n union

$O(n)$ makeSet

$O(m)$ find

$O(m)$ Union

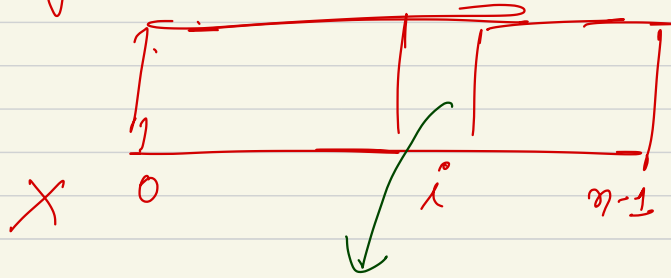
You can implement
Union-find data
structures in any
way & use here
→ but you also
want efficiency

Implementation view:

Simplest implementation

n people

Use an array?



$X[i]$ stores the identity of the group containing i th element

find(x): $O(1)$

↳ Simply look at the index of array

Define a group-id to each element

Ex. 3 7 11 12

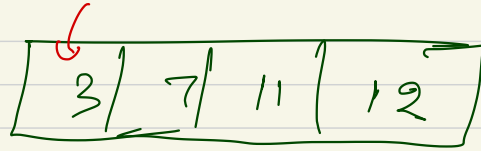
X

3	7	11	12
---	---	----	----

makeSet
find()
Union

Union

3 7 11 12



↑ ↑

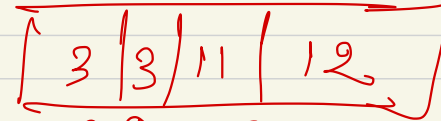
(3, 7)

m. Find → $O(m)$

n. Union → $O(n^2)$

Union(x, y)

Find (7) = 3



↑ ↑ ↑

↑

↑

①

②

③

$O(m+n^2)$
 Not very efficient

Change the identity of each element in B to A

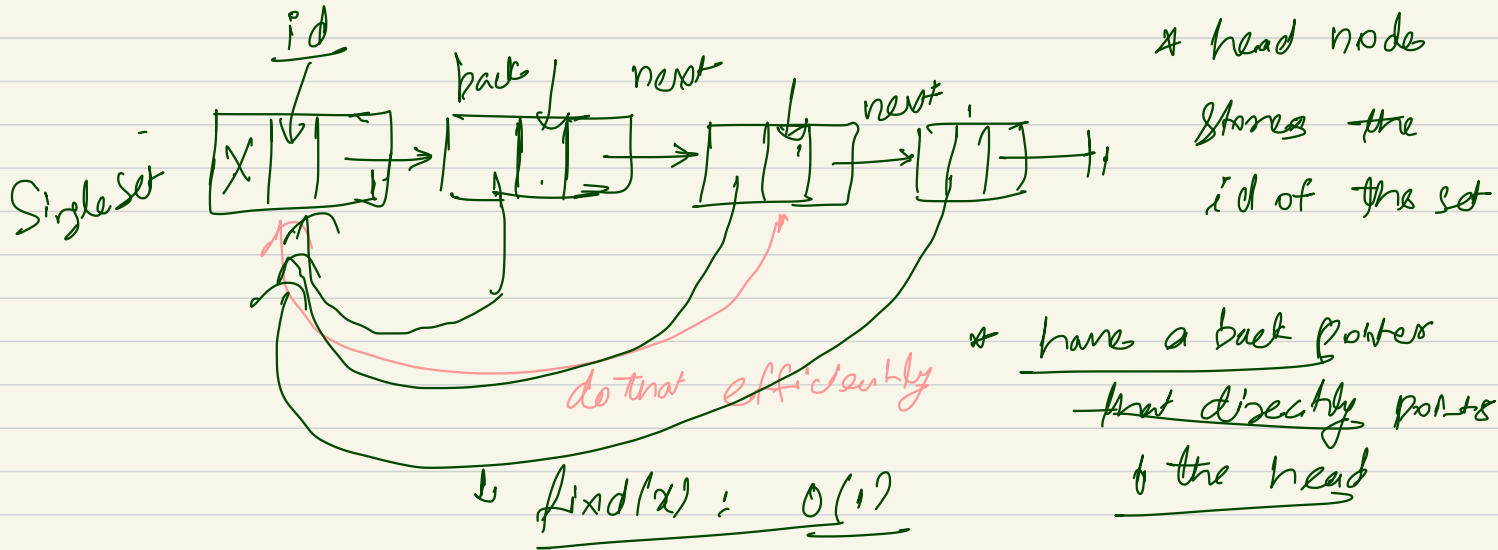
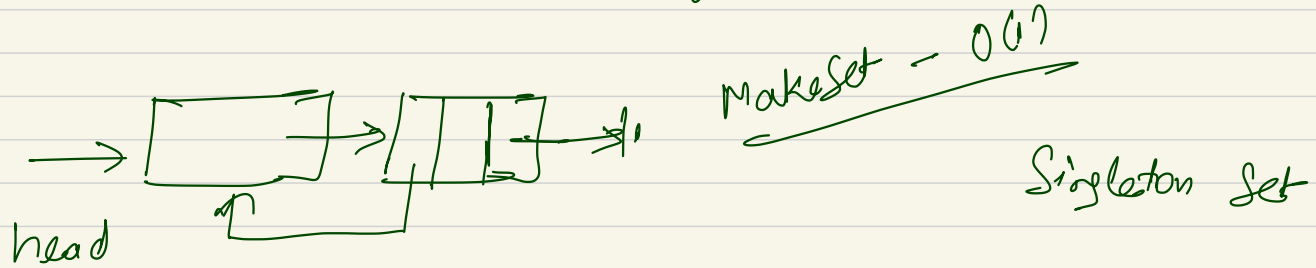
Set containing y → Set containing x

Complexity: worst case: → $O(n)$ complexity

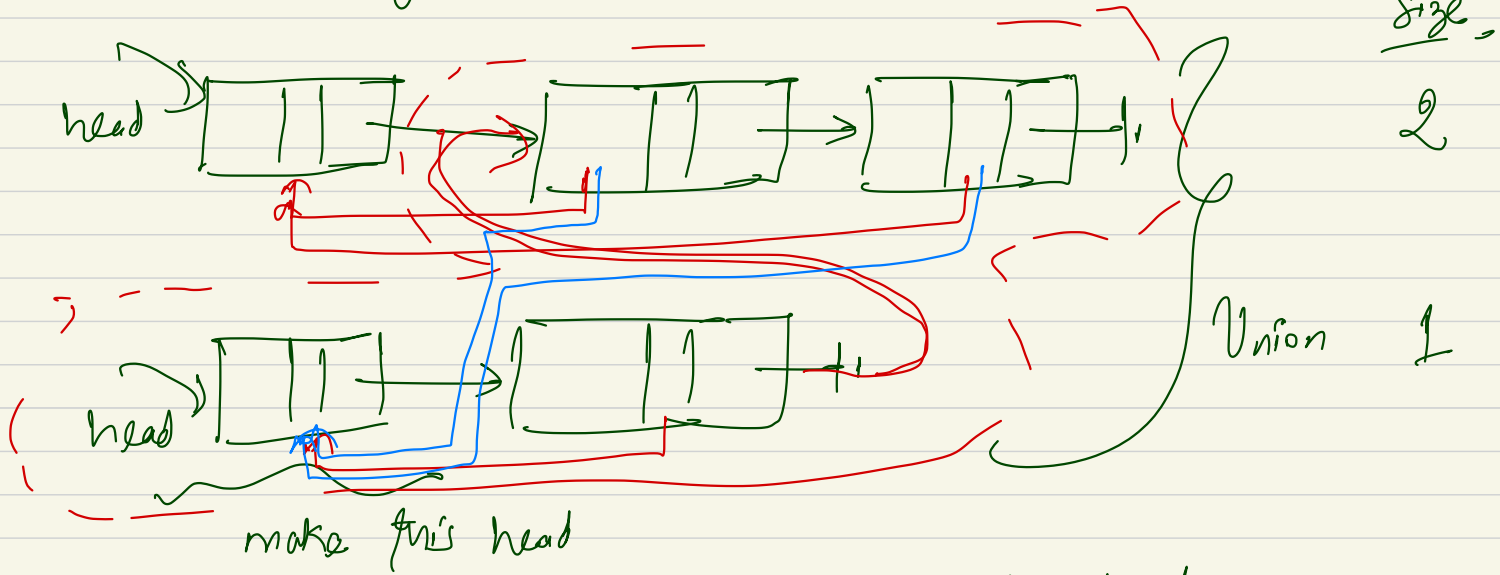
$b_1 \dots b_m$ A

Linked List Data Structure

Each set is denoted by a linked list



Union (x, y)



* add one list to the tail of another (B)
(single pointer update)

Worst case

$O(m)$

← *

Change back pointers of all the nodes in list A.

⇒

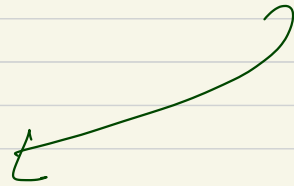
$O(m+n^2)$?? Inefficient

* head & tail pointers of each list

Can we do better?

1. Keep # of nodes in each list (in head)
2. During union, add the smaller list to end of the larger list

— will that help?

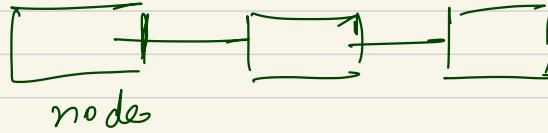


Worst Case - Smaller list
Can be $n/2$ etc

→ Changing $\frac{n}{2}$ back-pointers

Let's look at number of changes we're doing.

Look at a node, find an upper bound on the total # of times its back-pointer is updated.



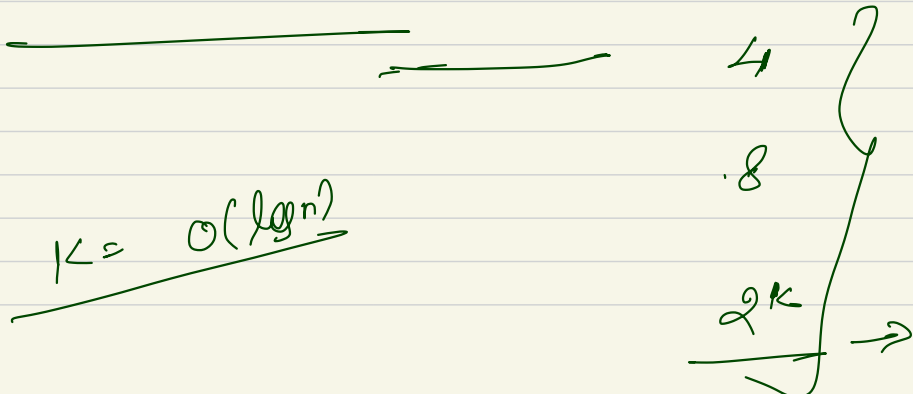
Change its backpointer - when it's appended at the end of a larger list

First time its back pointer was updated \rightarrow size of the resultant linked list \geq

2nd time

3rd

kth time



Single node: atmost $O(\log n)$ updates

n nodes \Rightarrow $O(n \log n)$ updates

m find + n Union

$O(m + n \log n)$

Linked list

$O(m + n^2)$

Static Array

Worst Case: Each operation adding a list to the tail of another list
 $\rightarrow O(n)$ time

In a series of operations, the worst case complexity for any one operation may be high, but on an average, complexity

be low. This average over a sequence gives a more reasonable estimate.

Amortized Analysis : Independently worst case may be high

Sequence of operations: average analysis

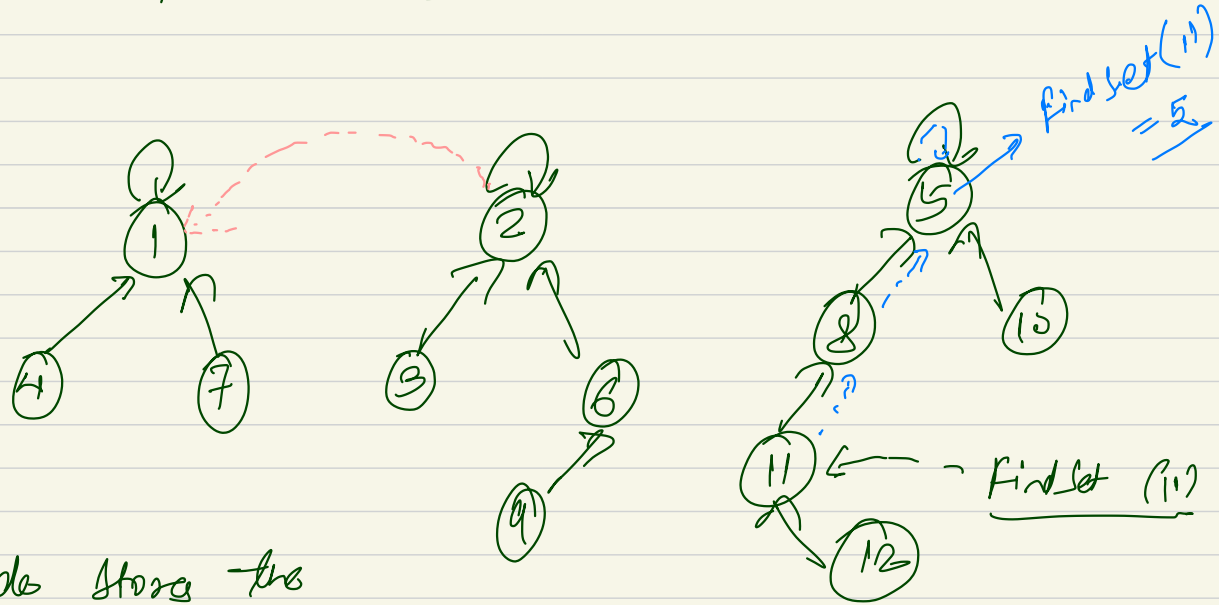
for a node $\rightarrow O(\log n)$ amortized complexity for back-pointer update

$O(n \log n)$

Another popular implementation - via rooted trees.

→ Each set is a rooted tree (inverted tree)
Id of a set → Id of root

3 sets



- * Each node stores the element of the set, and a parent pointer to the parent node.
- * The parent pointer of the root points to itself.

MakeSet (x) \Rightarrow



$O(1)$

find Set (x) \Rightarrow follow parent pointers to the root

Union (x, y) \Rightarrow Make root of one tree as child of another.

\downarrow
if you've the set IDs already

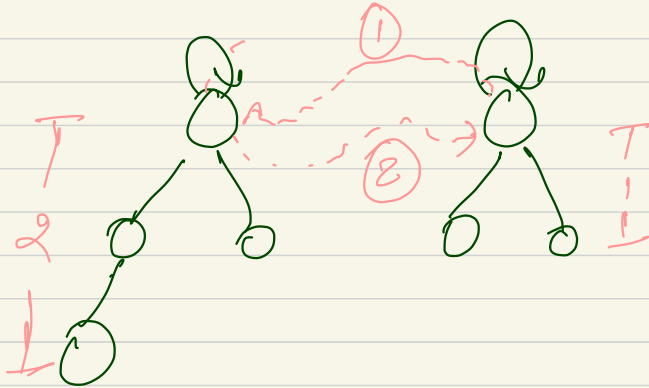
$O(1)$ ops

Worst case: this can be $O(n)$

Complexity
 $O(m+n)$ not efficient

Can I do better (balancing)?

Linked list



add a set A
to the tail of set B

set A is smaller

Follow 1: height = 2

Follow 2: height = 3

⇒ When a union is performed, the pointer of the smaller group is set to point to root of the larger group.

If we use this balancing, any tree of height h must have at least 2^h elements.

Proof \Rightarrow Indⁿ

Base Case: height 1



2 elements

Union

A B
(larger)

resultant trees

$$h = \max\{h(A), h(B) + 1\}$$

Induction step: Assume that both A & B follow the property $\geq 2^{h(A)} \hookrightarrow \geq 2^{h(B)}$

Case 1: $h(A)$ is larger

$$h = \max(h(A), \underline{h(B)+1})$$

\Rightarrow The final set has the height as that of A

but has more elements

$$n(A) \geq 2^{h(A)}$$

$$n(\text{new}) \geq n(A)$$

$$\geq 2^{h(A)}$$

$$= 2^{h(\text{new})}$$

Case 2: $h(B)+1$ is larger

$$n(B) \geq 2^{h(B)}$$

$$\underline{n(\text{new})} \geq \underline{2 \cdot n(B)}$$

$$\geq 2^{h(B)+1} = \underline{2^{h(\text{new})}}$$

\rightarrow both cases: new tree
will contain at least
 2^h elements

\hookrightarrow Proof

⇒ Using balancing, ensure trees have at least \log^h element
size $O(\log n)$ n elements

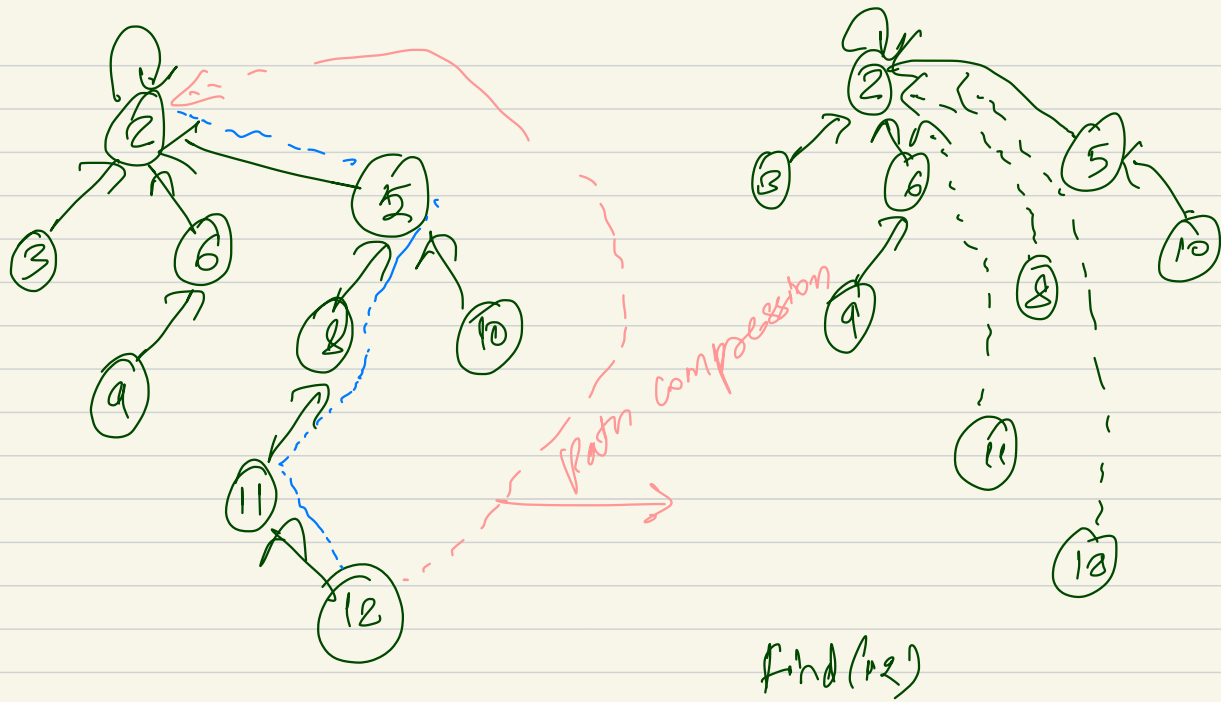
→ find operation can be done in $O(\log n)$ time

$$O(n + m \cdot \log n)$$

quite similar to linked list

We can do better than this?

Idea Path Compression ⇒ In a find operation, for each node v , that find visits, reset the parent pointer from v to point to the root.



Analysis shows that using path compression for a sequence of m operations of $\text{find} \geq n$ has a complexity

$$O(m \log^* n)$$

$\log^* n$ is recursively defined as

$$\log^* 1 = \log^* 2 = 1$$

for $n > 2$

$$\log^* n = 1 + \log^* (\lfloor \log_2 n \rfloor)$$

$$\text{Ex } \log^* 4 = 1 + \log^* 2 = 2$$

$$\log^* 14 = 1 + \log^* 4 = 3$$

$$\log^* \underline{60000} = 1 + \log^* 16 = 4$$

$\underline{< 2^{16}}$

for any number $n \leq \underline{2^{65536}}$

$\boxed{\log^* n \leq 5}$ \rightarrow as good as constant
for all practical purposes

⇒ Path compression helps in achieving

$O(m \log^* n + n)$ complexity

m finds

n Unions

⇒ $O(m \log^* n)$ for any practical input.

makeSet()

for each singleton element, x do

$x \cdot \text{parent} = x$

$x \cdot \text{size} = 1$

} // For union by size,
the element x stores
the size.

Algo Union (x, y) // Assume x & y are individual roots

if $x.size < y.size$

$x.parent = y$

$y.size = y.size + x.size$

$a = find(x)$
 $b = find(y)$
Union(a, b)

else

$y.parent = x$

$x.size = y.size + x.size$

Algo

find(x)

$x = x$

while $x.parent \neq x$

$x = x.parent$ // find the root

$z = x$

while $z.parent \neq z$

$w = z$

$z = z.parent$

$w.parent = x$ //

path

compression

ADT

Disjoint Sets / Union find

MakeSet(x)

Find(x)

Union(x,y)

Operⁿ View

Impⁿ View

Array

Linked List

Trees

Amortized Analysis

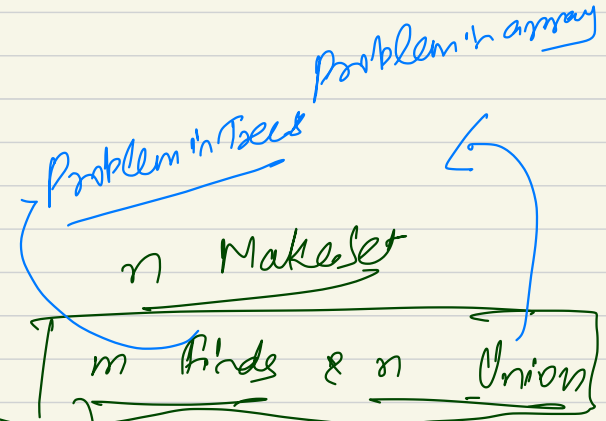
Sequence of n operations, worst

For an algo that needs

Case appears once in a while but

Kruskal's algo

any analysis gives a more realistic picture.



↓
Tree based implementation

$O(m+n)$ is the best