

Data Structures

Disjoint Sets 2

CS 225
Brad Solomon

March 25, 2026



UNIVERSITY OF
ILLINOIS
URBANA - CHAMPAIGN

Department of Computer Science

Learning Objectives

Continue to improve implementation of disjoint sets

Discuss how improvements affect efficiency

Find

$O(1)$



$O(\log n)$



$O(1)^*$

Union

$O(n)$

Smart
union

$O(\log n)$

Path
compression

$O(1)^*$

[unproven!]

(new proof)

Disjoint Sets

ADT:

`makeSet(vector<T> items)`

`Find(T key)`

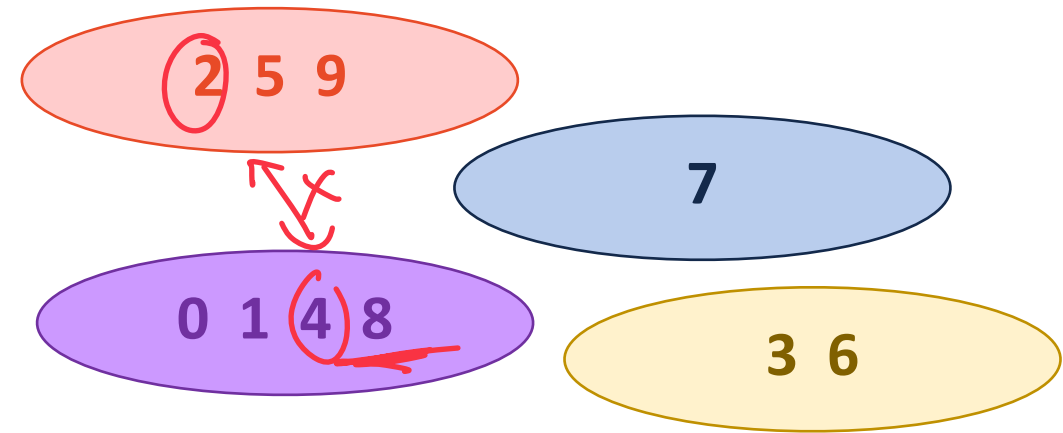
`Union(T k1, T k2)`

Key Ideas:

Every item exists in exactly one set

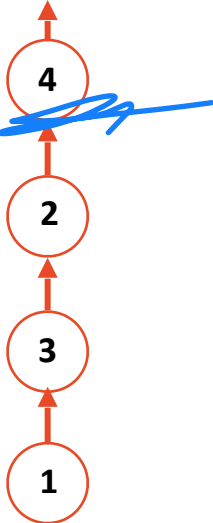
Every item in each set has same representation

Every set has a different representation



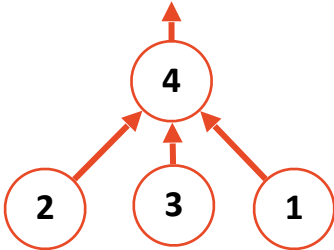
Disjoint Sets – Best and Worst UpTree

Reminder (values as index) can be anything!



PN 6 X → assign value ↑

100 million → 5
 -100 is key → 5



Node access by value as index

0	1	2	3	4
	3	4	2	-1

Stores 'Pointer' (or index) of parent node

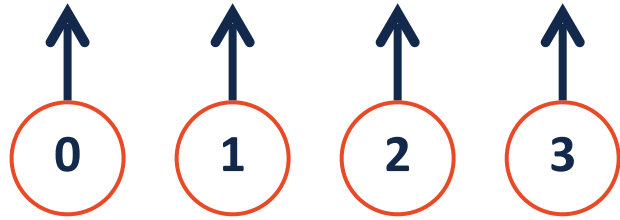
Key/Value we are looking for

0	1	2	3	4
	4	4	4	-1

Disjoint Sets Union by Size



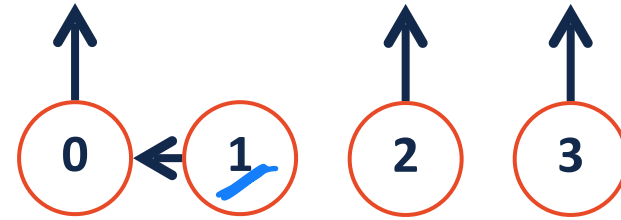
Join Code: 225



Store -size!

A

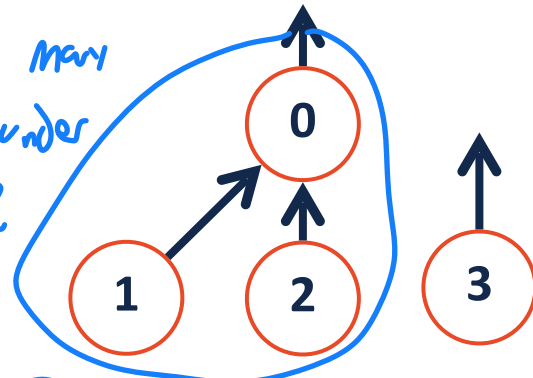
0	1	2	3
-1	-1	-1	-1



B

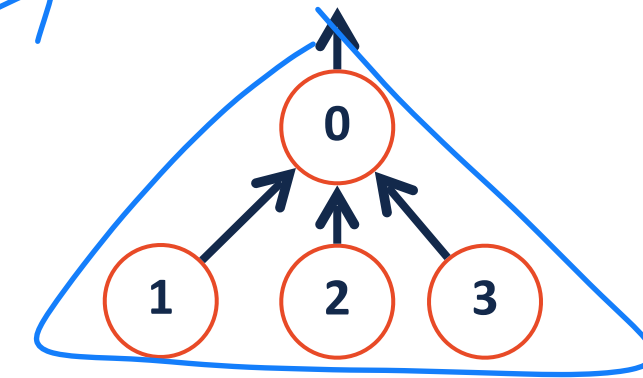
0	1	2	3
-2	0	-1	-1

How many nodes under root?



C

0	1	2	3
-3	0	0	-1



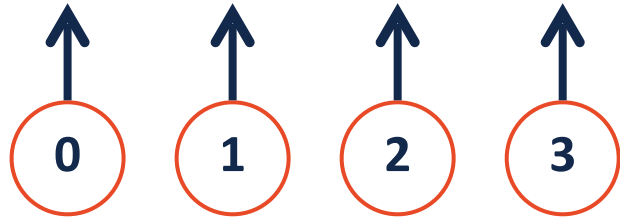
D

0	1	2	3
-4	0	0	0

Disjoint Sets Union by Size

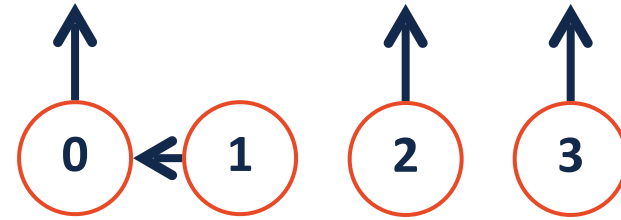


Join Code: 225



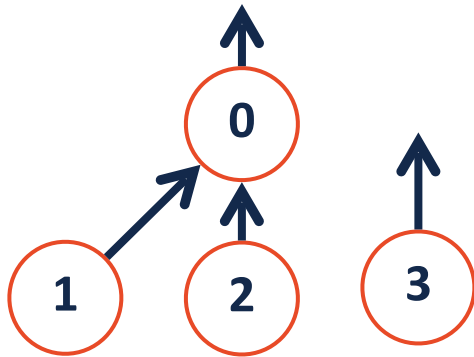
A

0	1	2	3
-1	-1	-1	-1



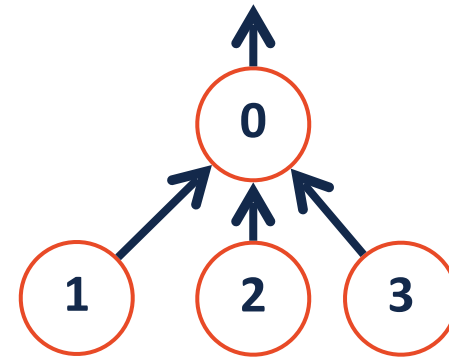
B

0	1	2	3
-2	0	-1	-1



C

0	1	2	3
-3	0	0	-1



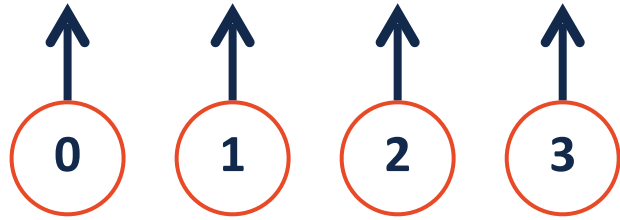
D

0	1	2	3
-4	0	0	0

Disjoint Sets Union by Height

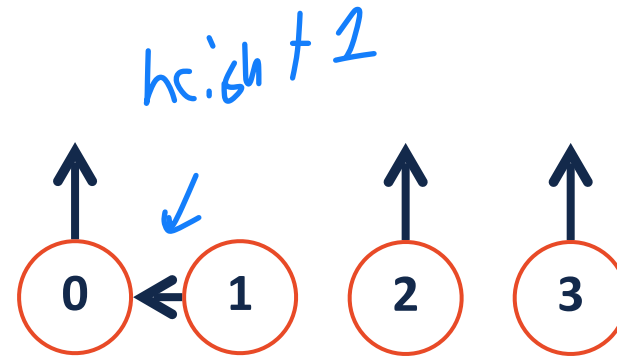


Join Code: 225



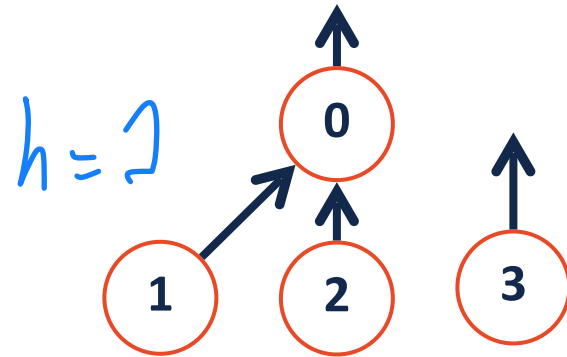
A

0	1	2	3
-1	-1	-1	-1



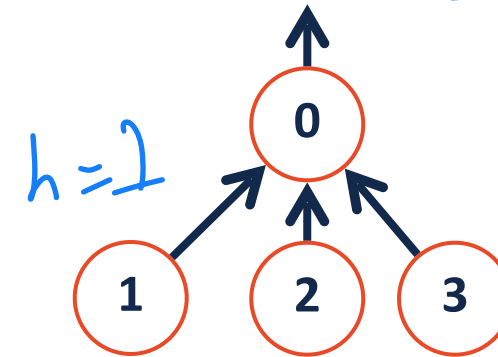
B

0	1	2	3
-2	0	-1	-1



C

0	1	2	3
-2	0	0	-1



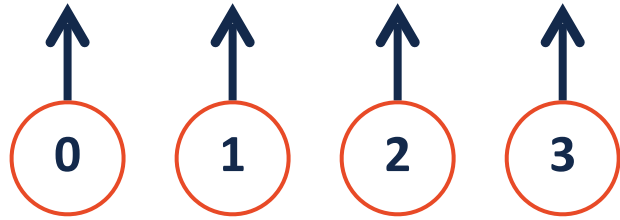
D

0	1	2	3
-2	0	0	0

Disjoint Sets Union by Height

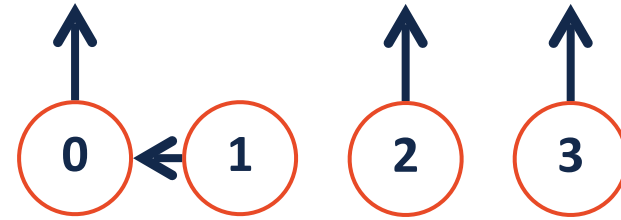


Join Code: 225



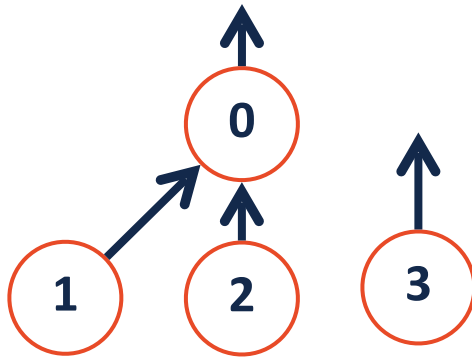
A

0	1	2	3
-1	-1	-1	-1



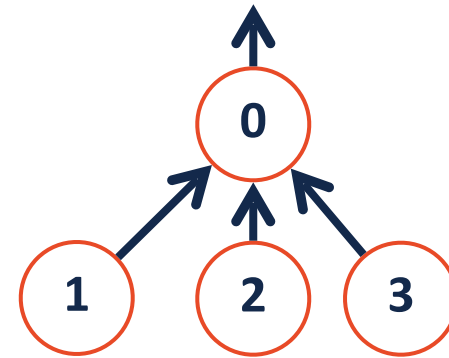
B

0	1	2	3
-2	0	-1	-1



C

0	1	2	3
-2	0	0	-1



D

0	1	2	3
-2	0	0	0

Disjoint Sets – Smart Union

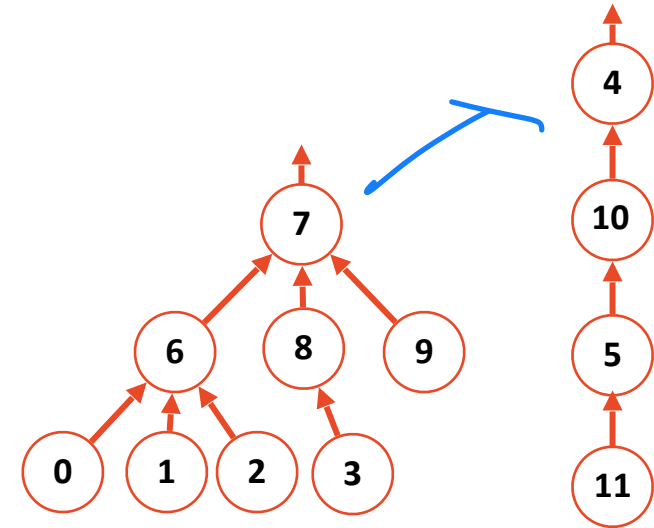


Two $O(1)$ methods of combining two sets

Claim: Both limit height to: $O(\log n)$.



this increases in height!



Union by height

Before Union

After Union

4	...	7
-4		-3

4	...	7
-4		4

Union by size

4	...	7
-4		-8

4	...	7
7		-12

Idea: Keep the height of the tree as small as possible.

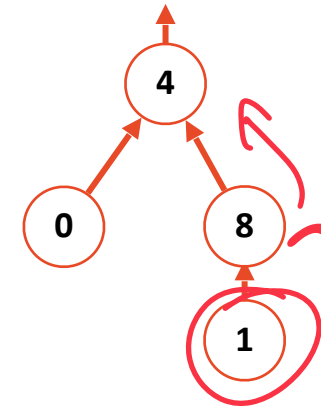
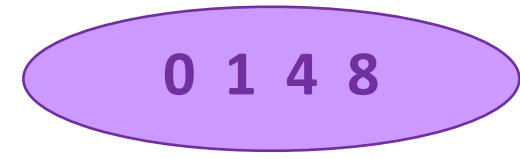
Idea: Minimize the number of nodes that increase in height

Disjoint Sets Find

Find(1)

```
1 int DisjointSets::find(int i) {  
2     if ( s[i] < 0 ) { return i; }  
3     else { return find( s[i] ); }  
4 }
```

*if negative
(canonical)*



Does implementation work on **height / size**?

Yes our find remains unchanged!

0	1	2	3	4	5	6	7	8	9
4	8			-3/-4				4	



Disjoint Sets Union

(cancel out root!) **unionBySize(4, 3)**

```

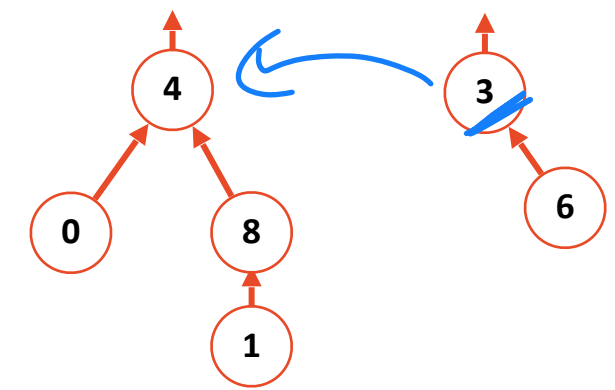
1 void DisjointSets::unionBySize(int root1, int root2) {
2   int newSize = arr[root1] + arr[root2]; O(1)
3
4   if ( arr[root1] < arr[root2] ) {
5     arr[root2] = root1; O(1)
6
7     arr[root1] = newSize; O(1)
8
9   } else {
10
11     arr[root1] = root2; O(1)
12
13     arr[root2] = newSize; O(1)
14
15   }
16 }

```

Secret step 0:
call find on orig large

Lookup both
update size!

Bigger absolute value



root2
root1

Don't forget those old neg values!

0	1	2	3	4	5	6	7	8	9
4	8		-2	-4		3		4	

-size -size

Disjoint Sets Union by Size

Claim: Sets unioned by size have a height of at most $O(\log_2 n)$

Claim: An UpTree of height h has nodes $\geq \frac{2^h}{n}$

↳ If true Think proof by induction!

Base Case:

$$h=0$$



1 node

$$1 \geq 1$$

$$1 \geq 2^0$$

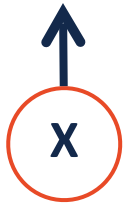
1 ✓

Disjoint Sets Union by Size

Claim: Sets unioned by size have a height of at most $O(\log_2 n)$

Claim: An UpTree of height h has nodes $\geq 2^h$

Base Case: $h = 0$



Base case height is 0, has one node.

vs.

$$2^0 = 1$$



Base case holds!

Disjoint Sets Union by Size

Claim: An UpTree of height h has nodes $\geq 2^h$

IH:

Disjoint Sets Union by Size

Claim: An UpTree of height h has nodes $\geq 2^h$

IH: Claim is true for $< i$ unions, prove for i th union (sets A and B).

(We have done $i - 1$ total unions and plan to do **one** more)

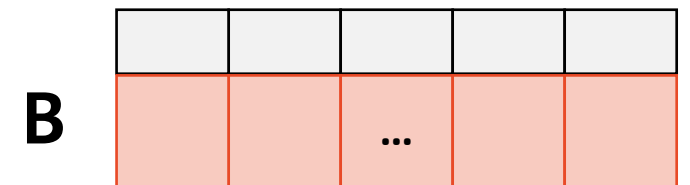
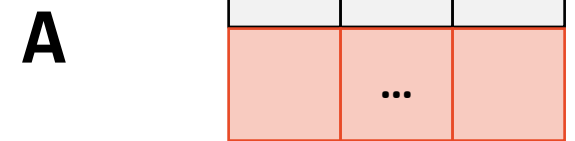
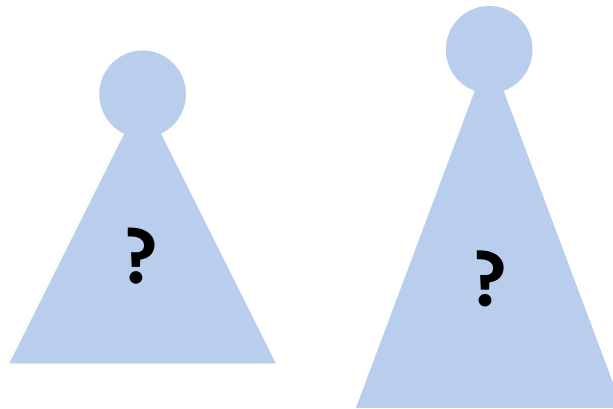
Without loss of generality, let B be the larger set **BY SIZE**

We must explore how height changes for each case:

Case 1: $h(A) < h(B)$

Case 2: $h(A) == h(B)$

Case 3: $h(A) > h(B)$



size(B) \geq size(A)

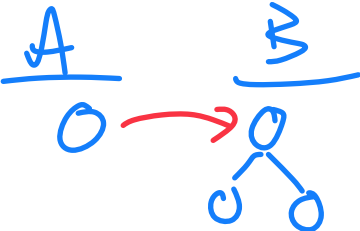
$$\text{size}(B) \geq \text{size}(A)$$

Disjoint Sets Union by Size

Claim: An UpTree of height h has nodes $\geq 2^h$

IH: Claim is true for $< i$ unions, prove for i th union (sets A and B).

Case 1: height(A) < height(B)



100% true here:

Unioning A into B does not increase height!

Union smaller size into larger,

B' is new upTree $\rightarrow h(B') = h(B)$

$$\text{size}(B) \geq \text{size}(A)$$

Disjoint Sets Union by Size

Claim: An UpTree of height h has nodes $\geq 2^h$

IH: Claim is true for $< i$ unions, prove for i th union (sets A and B).

Case 1: height(A) < height(B)

Ideal case where size and height in agreement!

Height doesn't change ($h(B') = h(B)$).

By IH: $\boxed{\text{size}(A) \geq 2^{h(A)} \quad \text{size}(B) \geq 2^{h(B)}}$

$$\boxed{\text{size}(B')} = \text{size}(A) + \text{size}(B) = 2^{h(A)} + 2^{h(B)} \geq 2^{h(B)} = \boxed{2^{h(B')}}$$

By IH $(x+1) \geq x$
At least 1



$$\text{size}(B) \geq \text{size}(A)$$

Disjoint Sets Union by Size

Claim: An UpTree of height h has nodes $\geq 2^h$

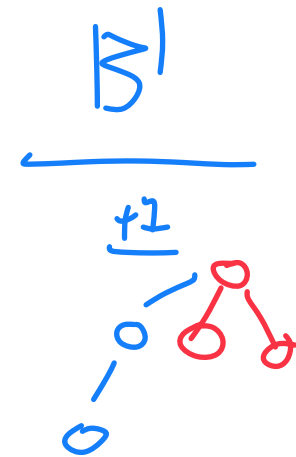
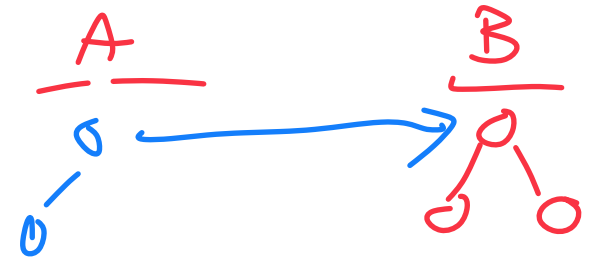
IH: Claim is true for $< i$ unions, prove for i th union (sets A and B).

Case 2: $\text{height}(A) == \text{height}(B)$

$$h(B') = h(B) + 1$$

or

$$h(A) + 1$$



$$\text{size}(B) \geq \text{size}(A)$$

Disjoint Sets Union by Size

Claim: An UpTree of height h has nodes $\geq 2^h$

IH: Claim is true for $< i$ unions, prove for i th union (sets A and B).

Case 2: $\text{height}(A) == \text{height}(B)$

$$h(B') = h(B) + 1$$

If we merge two equal height trees, height always increase by 1

By IH: $\left[\text{size}(A) \geq 2^{h(A)} \quad \text{size}(B) \geq 2^{h(B)} \right]$

$$\text{size}(B') = \text{size}(A) + \text{size}(B) = 2^{h(A)} + 2^{h(B)} \quad \text{w/c } h(A) = h(B)$$

$$= 2^{h(B)} + 2^{h(B)}$$

$$= 2 * 2^{h(B)} = 2^{h(B)+1} \geq 2^{h(B')}$$

Disjoint Sets Union by Size

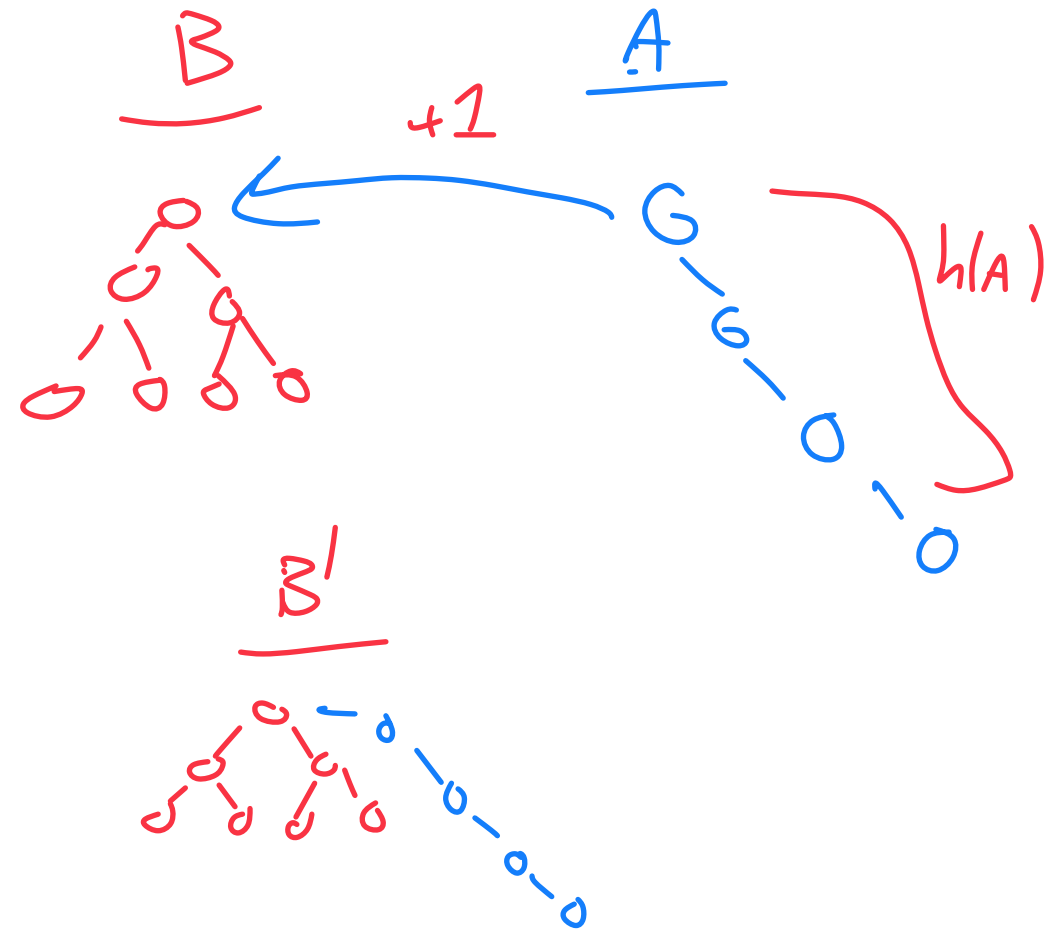
$$\underline{\text{size}(B)} \geq \text{size}(A)$$

Claim: An UpTree of height h has nodes $\geq 2^h$

IH: Claim is true for $< i$ unions, prove for i th union (sets A and B).

Case 3: height(A) > height(B)

$$h(B') = h(A) + 1$$



Disjoint Sets Union by Size

$$\underline{\text{size}(B)} \geq \underline{\text{size}(A)}$$

Claim: An UpTree of height h has nodes $\geq 2^h$

IH: Claim is true for $< i$ unions, prove for i th union (sets A and B).

Case 3: height(A) > height(B)

$$h(B') = h(A) + 1$$

Merging taller tree into smaller — height increase to height(A)+1!

By IH: $\text{size}(A) \geq 2^{h(A)}$ $\text{size}(B) \geq 2^{h(B)}$

$$\text{size}(B') = \text{size}(A) + \text{size}(B) \geq 2 \text{size}(A)$$

$$\begin{aligned} &\downarrow \\ &\geq \text{size}(A) + \text{size}(A) \end{aligned}$$

$$= 2 * 2^{h(A)} = 2^{h(A)+1} \geq 2^{h(B')}$$

Disjoint Sets Union by Size

$$\text{size}(B) \geq \text{size}(A)$$



Proven: An UpTree of height h has nodes $\geq 2^h$

IH: Claim is true for $< i$ unions, prove for i th union.

Each case we saw we have $n \geq 2^h$.



$O(\log n)$ Find
↓

$O(\log n)$ union (arbitrary start value)

union is $O(1)$
if given
(arbitrary)

Disjoint Sets Find

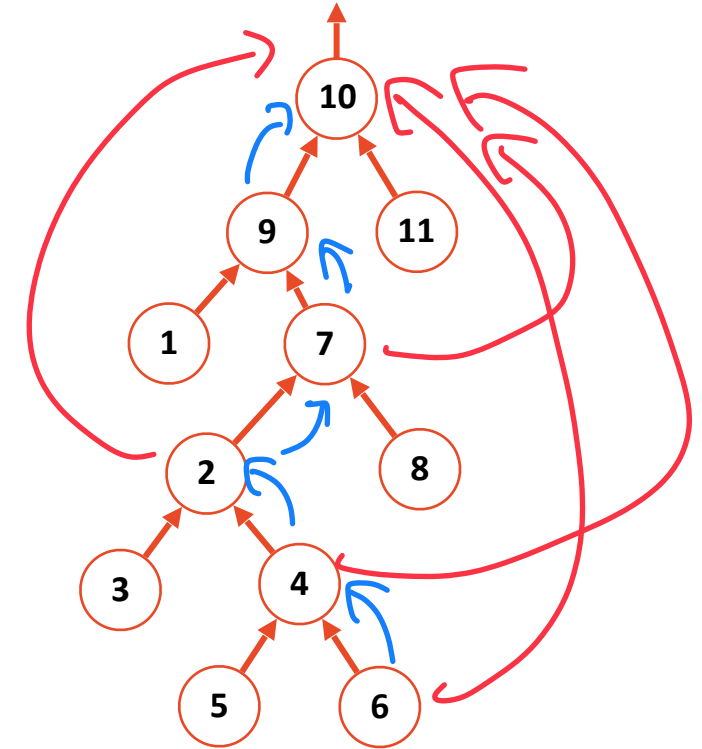
Find(6)

6 → 10

```
1 int DisjointSets::find(int i) {  
2     if ( s[i] < 0 ) { return i; }  
3     else { return find( s[i] ); }  
4 }
```

As we walk up a tree, why cant we fix it?

why not update my uptr?



Disjoint Sets Find

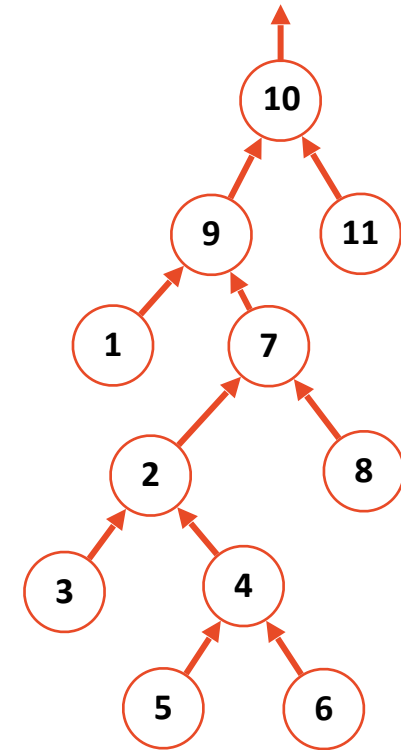
Find(6)

```
1 int DisjointSets::find(int i) {  
2     if ( s[i] < 0 ) { return i; }  
3     else { return find( s[i] ); }  
4 }
```

As we walk up a tree, why cant we fix it?

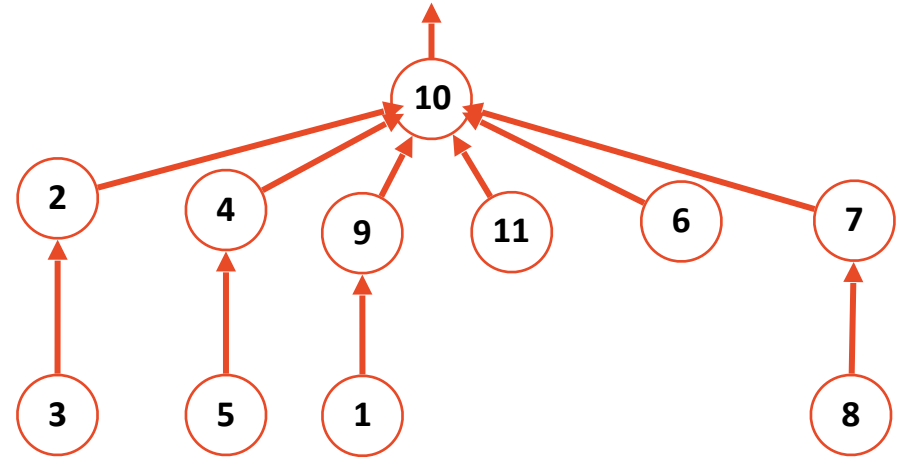
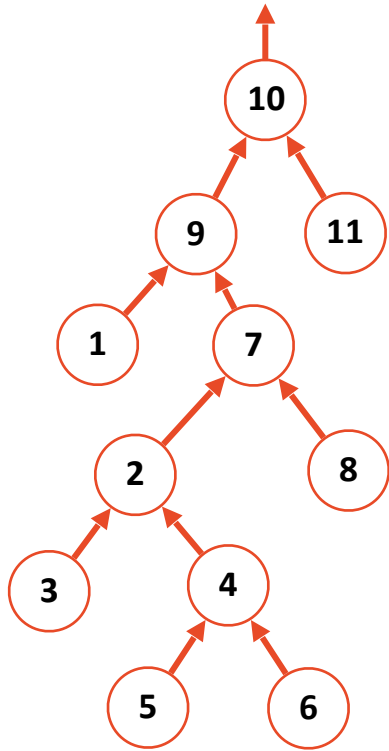
This is **path compression**:

```
1 int DisjointSets::find(int i) {  
2     if ( s[i] < 0 ) { return i; }  
3     else {  
4         int root = find( s[i] );  
5         s[i] = root;  
6         return root;  
7     }  
8 }
```



Path Compression

Find(6)



This seems good — but how good in theory?

$O(1)$ *

Path Compression Analysis

Two major problems here:

1) Our efficiency changes ***over repeated calls to find()***

2) Our height changes so we cant use union ~~by~~ height

Amortized Time Review

We have **n items**. We make **n insert()** calls.

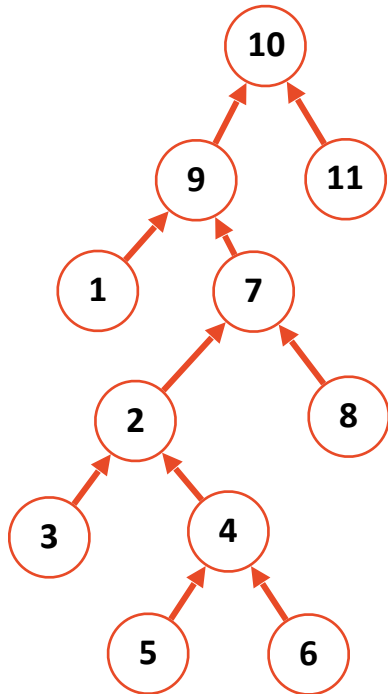
We are interested in the **worst case work** possible **over n calls**.



Amortized Time (Path Compression)

We have **n items** in an Uptree. We make **m find()** calls.

We are interested in the **worst case work** possible **over m calls**.



Union by Rank (Not Height)

Once I do path compression, I change the height of tree!

So we need a new way of approximating height.

Rank is a way of remembering what our height was before P.C.

Union by Rank (Not Height)

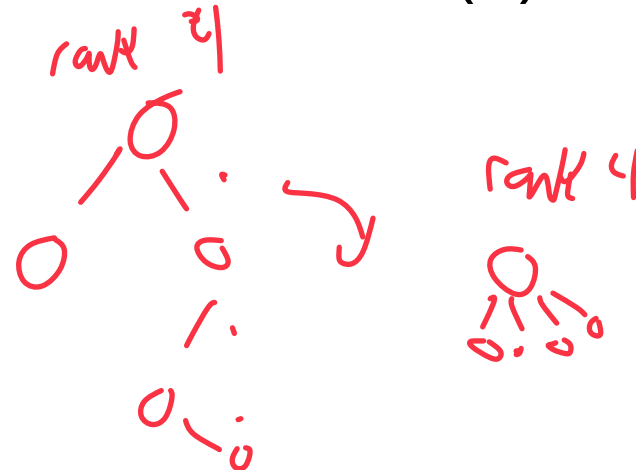
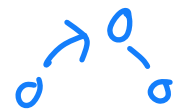
New UpTrees have rank = 0

Let A, B be two sets being unioned. If:

rank(A) == rank(B): The merged UpTree has rank + 1

rank(A) > rank(B): The merged UpTree has rank(A)

rank(B) > rank(A): The merged UpTree has rank(B)



Key Properties of UpTree by rank w/ PC

The parent of a node is always higher rank than the node.

There are at least $n \geq 2^r$ nodes in a root of rank r .

For any integer r , there are at most $\frac{n}{2^r}$ nodes of rank r .

Key Properties of UpTree by rank w/ PC

The parent of a node is always higher rank than the node.

This comes from how we set up rank union

(Take larger of two rank or add one if tied)

There are at least $n \geq 2^r$ nodes in a root of rank r .

Proof by Induction: To create rank r set, we merge two $r - 1$ sets

By IH (not shown), those sets have $2^{r-1} + 2^{r-1} = 2^r$ nodes

For any integer r , there are at most $\frac{n}{2^r}$ nodes of rank r .

A rewrite of the above logic given n nodes

Just no this is true

Amortized Time (Rank w/ Path Compression)

Put every non-root node in a bucket by rank!

Structure buckets to store ranks $[r, 2^r - 1]$

\uparrow \uparrow

Where did number range come from?

Ranks	Bucket
0	0
1	1
2 - 3	2
4 - 15	3
16 - 65535	4
65536 - $2^{\{65536\}} - 1$	5

\log^*

Size of universe
quadrillions

Iterated Logarithm Function ($\log^* n$)

The number of times you can take a log of a number

$$\log^*(n) = \begin{cases} 0 & , n \leq 1 \\ 1 + \log^*(\log(n)) & , n > 1 \end{cases}$$

$$\log^*(2^{65536}) = 5$$

This is ~~most~~ useful part of proof!

$$\begin{array}{l} \log 2^{65536} \\ \log 2^{16} = 65536 \\ \log 2^4 = 16 \\ \log 2^2 = 4 \\ 2^1 = 2 \\ 2^0 = 1 \end{array}$$

In CS, we consider this a constant
(But its not)

Amortized Time (Rank w/ Path Compression)

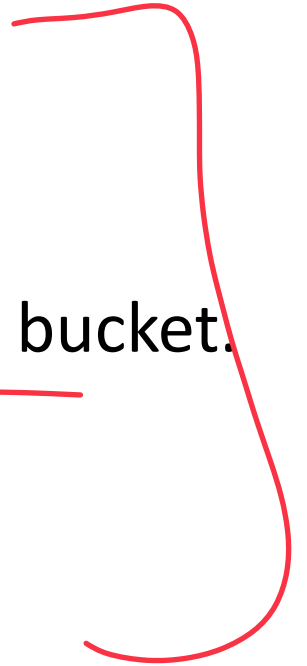
The work of **find(x)** are the steps taken on the path from a node x to the root (or immediate child of the root) of the UpTree containing x

We can split this into two cases:

Case 1: We take a step from one bucket to another bucket.

Case 2: We take a step from one item to another inside the same bucket.

Plus both
are $O(\log^* x)$



Amortized Time (Rank w/ Path Compression)

The work of **find(x)** are the steps taken on the path from a node x to the root (or immediate child of the root) of the UpTree containing x

We can split this into two cases:

Case 1: We take a step from one bucket to another bucket.

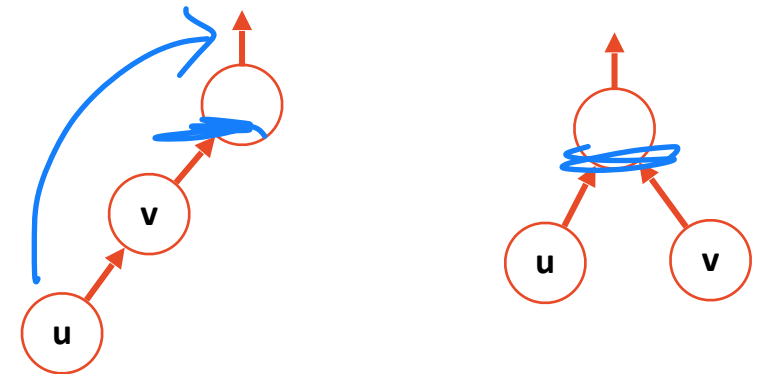
We have at most $\log^*(n)$ buckets so for m finds, this is $O(m \log^* n)$

Case 2: We take a step from one item to another inside the same bucket.

Let's call this the step from u to v .

Every time we do this, we do path compression:

We set $\text{parent}(u)$ a little closer to root



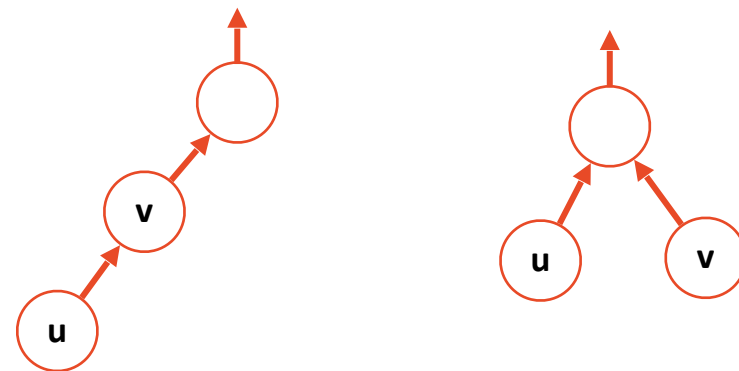
Amortized Time (Rank w/ Path Compression)

Case 2: We take a step from one item to another *inside* the same bucket.

Let's call this the step from **u** to **v**.

Every time we do this, we do path compression:

We set $\text{parent}(u)$ a little closer to root



How many total times can I do this for each **u** in a bucket?

By definition of our bucket ranges $\sim \frac{n}{2^r}$ \leftarrow # of times each can update

How many nodes are in bucket **r**?

By definition of how we set up rank: $\frac{n}{2^r}$ \leftarrow # of items

Given we have $\log^*(n)$ buckets:

Case 2 work is $n \log^*(n)$

Final Result



We have **n items** in an Uptree. We make **m find()** calls. Total work is:

Amortized $(n + m) \log^*(n)$

In terms of real world data, this is practically a constant.



We can only update paths so many times before call $u()$

↳ This is a pseudo-constant

Alternative Not-Actually-A-Proof

Unproven Claim: A disjoint set implemented with smart union and path compression with **m** find calls and **n** items has a worst case running time of **inverse Ackerman**. $[O(m \alpha(n))]$

This grows *very* slowly to the point of being treated a constant in CS.