

Data Structures

Disjoint Sets 2

CS 225

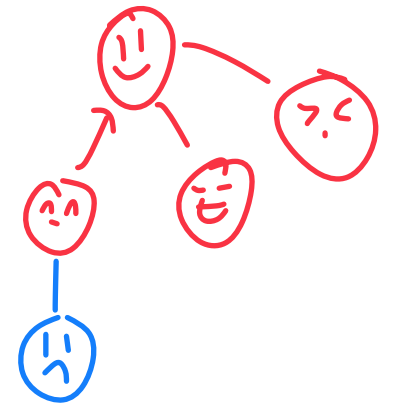
Brad Solomon

October 18, 2024



UNIVERSITY OF
ILLINOIS
URBANA - CHAMPAIGN

Department of Computer Science



Learning Objectives

Continue to improve implementation of disjoint sets

Discuss how improvements affect efficiency

Disjoint Sets

ADT:

`makeSet(vector<T> items)`

`Find(T key)`

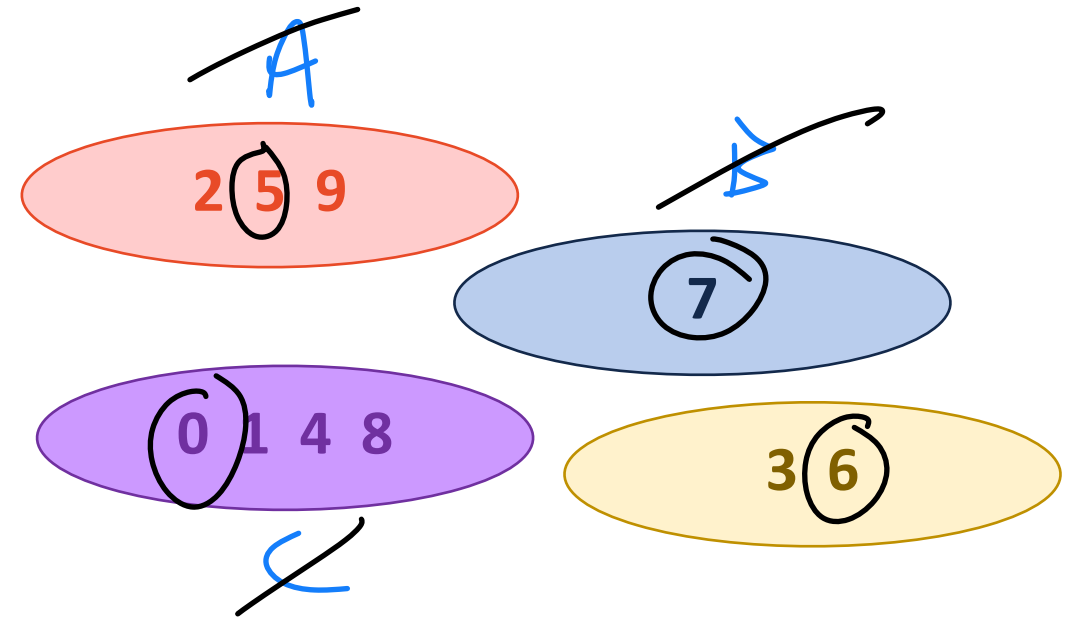
`Union(T k1, T k2)`

Key Ideas:

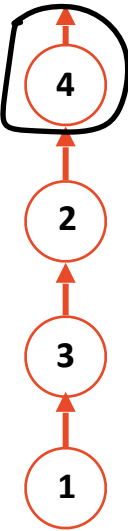
Every item exists in exactly one set

Every item in each set has same representation

Every set has a different representation

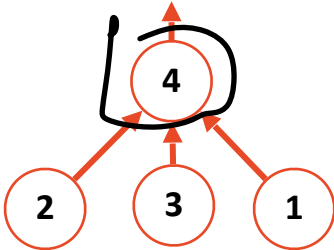


Disjoint Sets – Best and Worst UpTree



~~My Data Set~~
 My Data Set
 my PNG

 ↳ 5

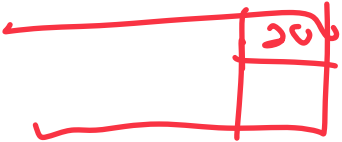


0	1	2	3	4
	3	4	2	-1

↑
 Blank!
 ↘

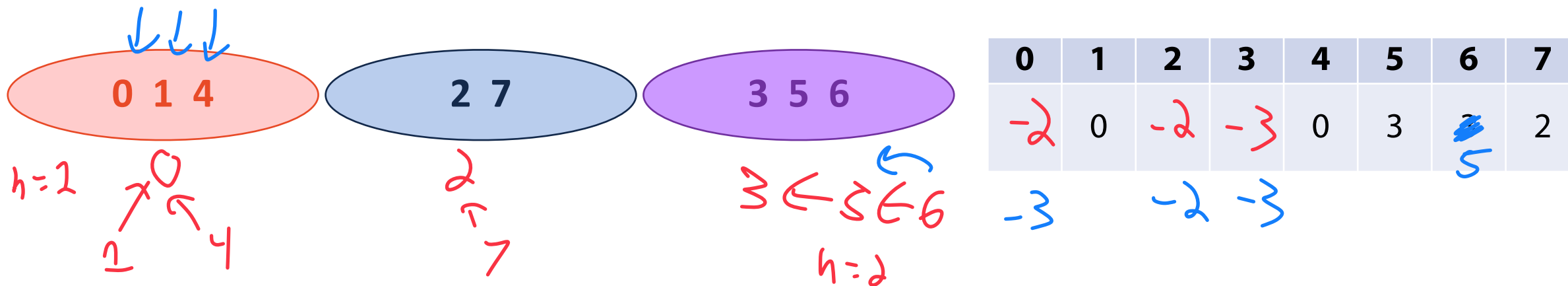
0	1	2	3	4
	4	4	4	-1

↑ ↘



Disjoint Set Implementation

Store an UpTree as an array, canonical items store **height** / **size**



Find(k): Repeatedly look up values until **negative value**

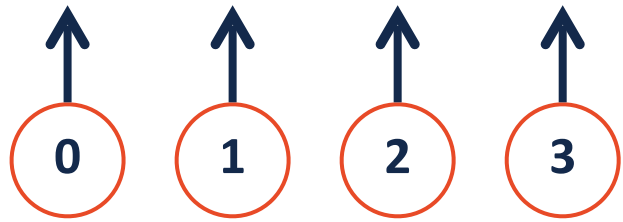
↳ Height: $-1 \cdot (\text{height} + 1)$ ② - value - 1

↳ Size: $-1 \cdot \text{size}$ ① - value - 1

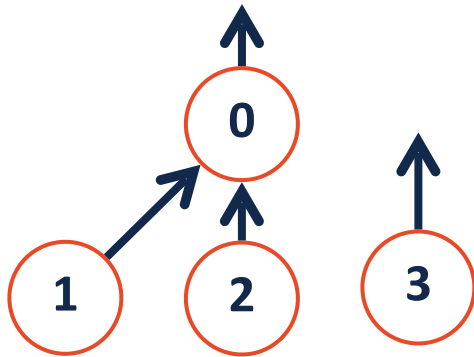
Union(k₁, k₂): Update **smaller** canonical item to point to larger

Update value of remaining canonical item

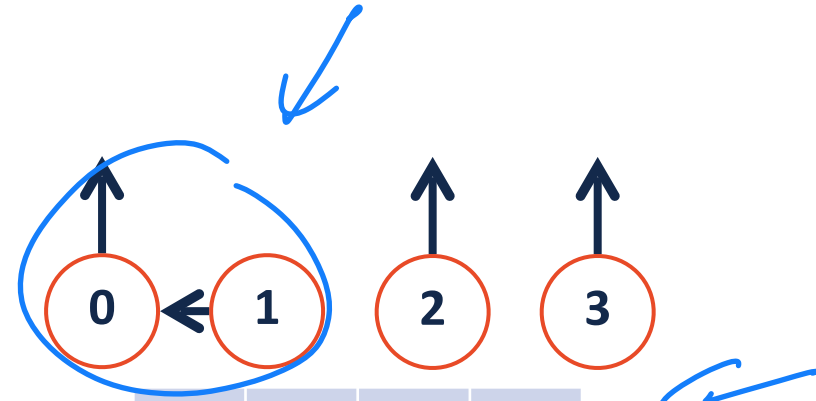
Disjoint Sets Union by Size



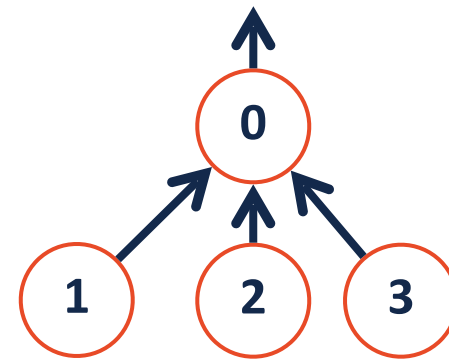
0	1	2	3
-1	-1	-1	-1



0	1	2	3
-3	0	0	-1



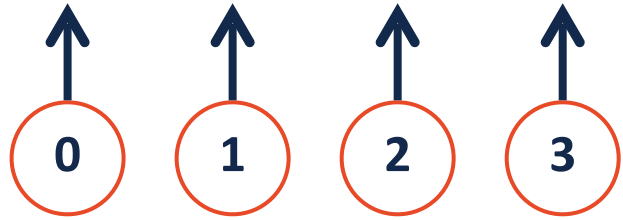
0	1	2	3
-2	0	-1	-1



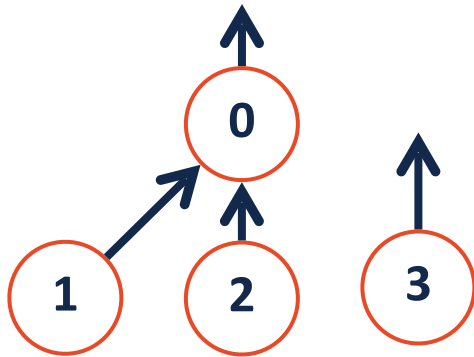
0	1	2	3
-4	0	0	0

!!!

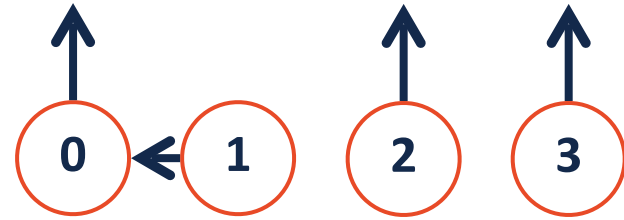
Disjoint Sets Union by Size



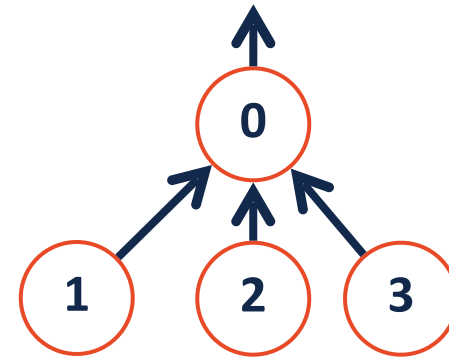
0	1	2	3
-1	-1	-1	-1



0	1	2	3
-3	0	0	-1

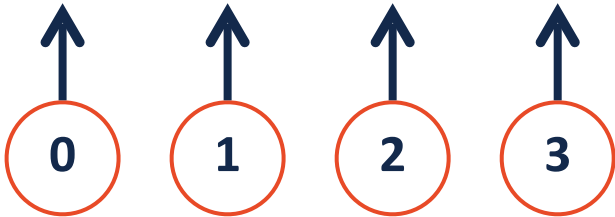


0	1	2	3
-2	0	-1	-1

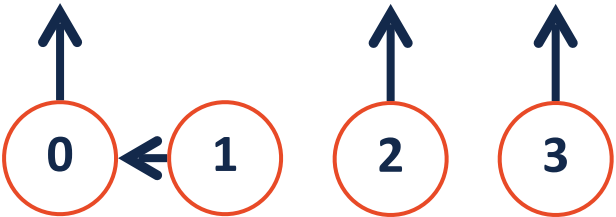


0	1	2	3
-4	0	0	0

Disjoint Sets Union by Height

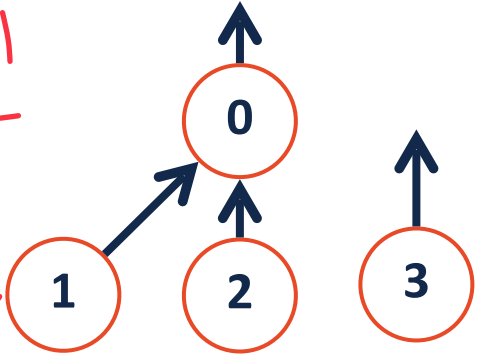


0	1	2	3
-1	-1	-1	-1

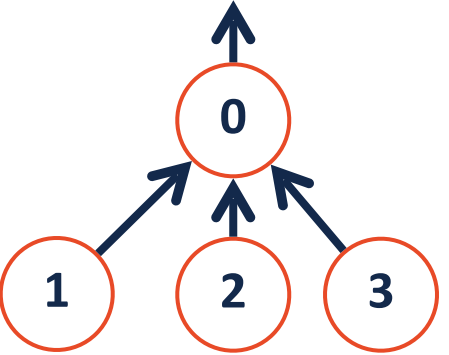


0	1	2	3
-2	0	-1	-1

$(h+1)_{a-1}$



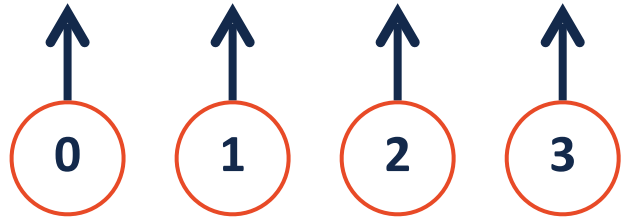
0	1	2	3
-2	0	0	-1



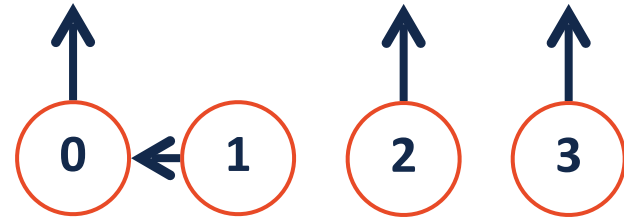
0	1	2	3
-2	0	0	0



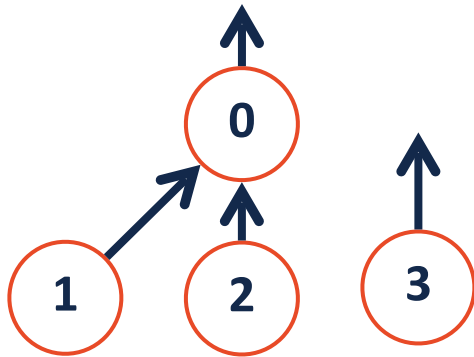
Disjoint Sets Union by Height



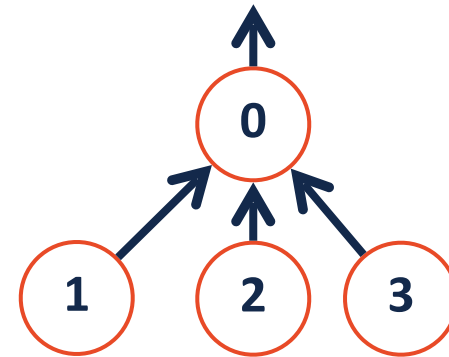
0	1	2	3
-1	-1	-1	-1



0	1	2	3
-2	0	-1	-1



0	1	2	3
-2	0	0	-1

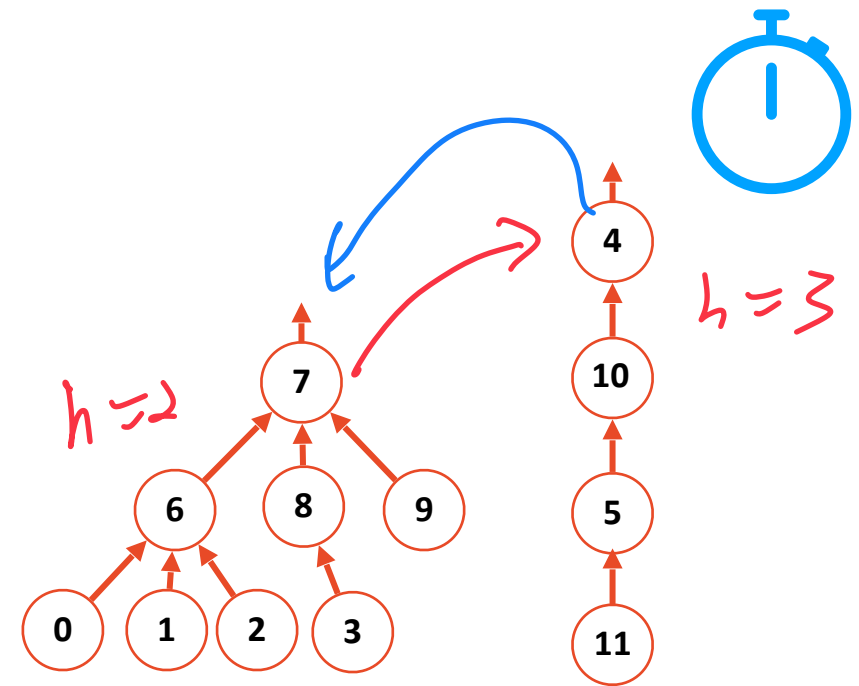


0	1	2	3
-2	0	0	0

Disjoint Sets – Smart Union

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

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



Union by height

Before Union

After Union

4	...	7
<u>-4</u>		<u>-3</u>

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

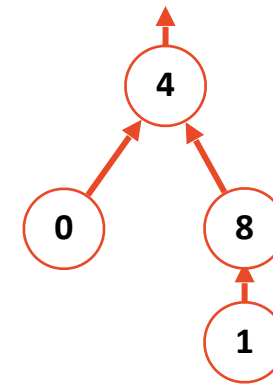
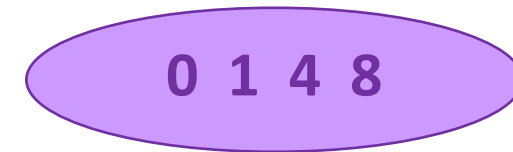
Find(1)

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

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

↳ Yes as long as canonical items
are negative

✓ why we set height
- (height + 1)



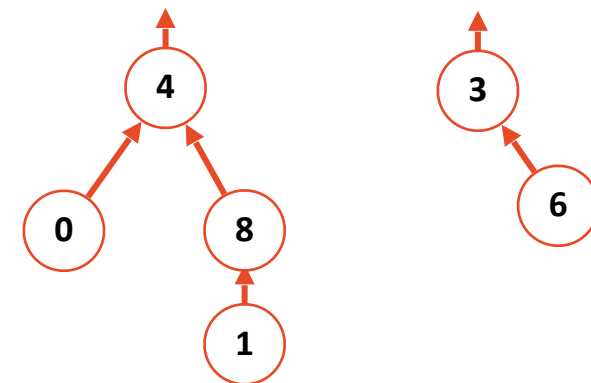
0	1	2	3	4	5	6	7	8	9
4	8			-3/-4				4	



Disjoint Sets Union

unionBySize(4, 3)

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



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



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 $\geq \underline{2^h}$

↳ proof by induction

Base Case: $n = 1$

$$h = 0$$

(A)

$$n \geq 2^0$$
$$1 \geq 1$$

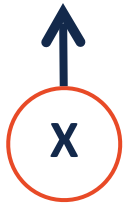


Disjoint Sets Union by Size

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: $h = 0$



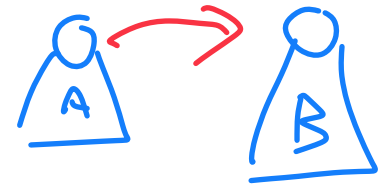
Base case height is 0, has one node.

vs.

$$2^0 = 1$$

Base case holds!

Disjoint Sets Union by Size



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

IH: Claim true for up to $i-1$. Prove for i .

Let A, B be two sets, Let B be the larger set. (Always union A into B)

Must show all cases of height

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

2. $h(A) = h(B)$

3. $h(A) > h(B)$

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)

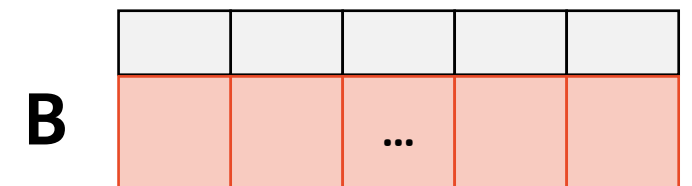
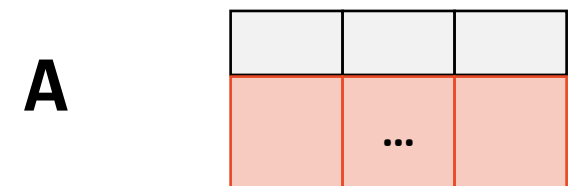
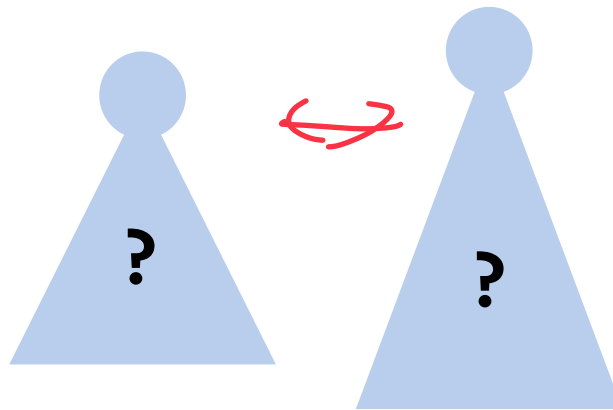
Without loss of generality, let B be the larger set **BY SIZE**

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)

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

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).

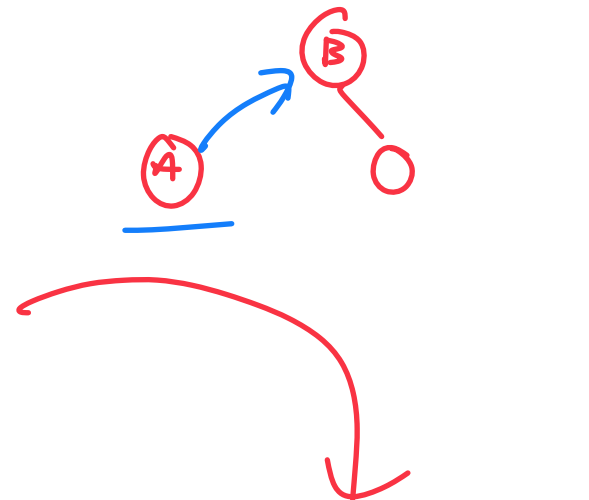
Case 1: height(A) < height(B)

Best case! Size + height agree.

Height doesn't change $\rightarrow \underline{h(B) = h(B')}$

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

$$\left[\begin{array}{l} \text{Size } (B') = 2^{h(A)} + 2^{h(B)} \\ \uparrow \text{new B} \quad \uparrow \text{By def} \end{array} \right] \geq 2^{h(B)} \geq 2^{h(B')} \quad \begin{array}{l} \uparrow \text{By logic} \end{array}$$



Disjoint Sets Union by Size

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

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

$$\wedge \geq 2^{h(B')}$$

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

Case 1: $\text{height}(A) < \text{height}(B)$

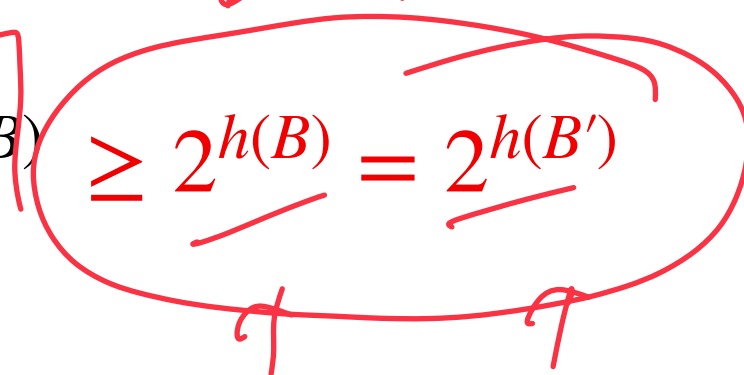
Ideal case where size and height in agreement!

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

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

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

\nearrow # items



Disjoint Sets Union by Size

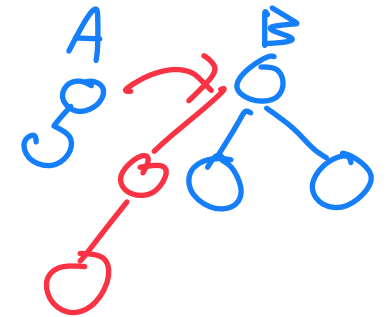
$$\text{size}(B) \geq \text{size}(A)$$

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

If merge two same height trees, height = height + 1



Disjoint Sets Union by Size

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

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

$$\begin{aligned} \text{size}(B') &= \text{size}(A) + \text{size}(B) = 2^{h(A)} + 2^{h(B)} \\ &= 2^{h(B)} + 2^{h(B)} \\ &= 2 * 2^{h(B)} = 2^{h(B)+1} \geq 2^{h(B')} \end{aligned}$$

Handwritten notes: A blue arrow points from $h(A) = h(B)$ to the first $2^{h(B)}$ term. A blue arrow points from the $2 * 2^{h(B)}$ term to $2^{h(B)+1}$. A red checkmark is at the end of the final inequality.

Disjoint Sets Union by Size

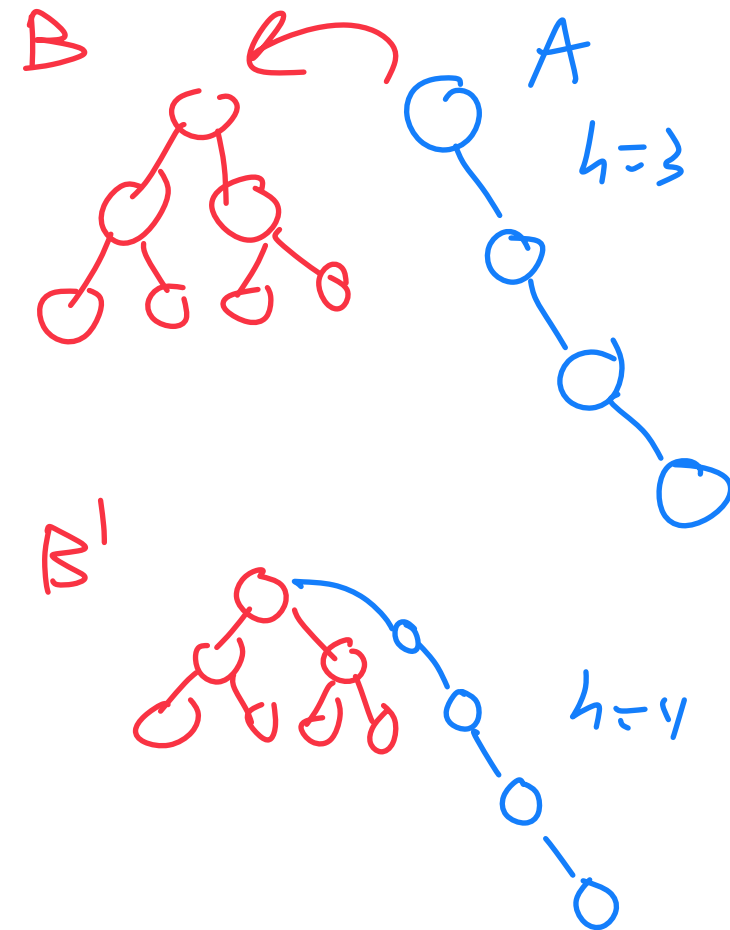
$$\text{size}(B) \geq \text{size}(A)$$

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: $\text{height}(A) > \text{height}(B)$

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



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

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 3: $\text{height}(A) > \text{height}(B)$

$h(B')$

Merging taller tree into smaller — **height increase to $\text{height}(A)+1$!**

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

$$\begin{aligned} \text{size}(B') &= \text{size}(A) + \text{size}(B) \geq 2 \text{size}(A) \\ &= 2 * 2^{h(A)} = 2^{h(A)+1} \geq 2^{h(B')} \end{aligned}$$

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 $n \geq 2^h$. for a set of height h

$$h \sim O(\log n)$$

Disjoint Sets Find

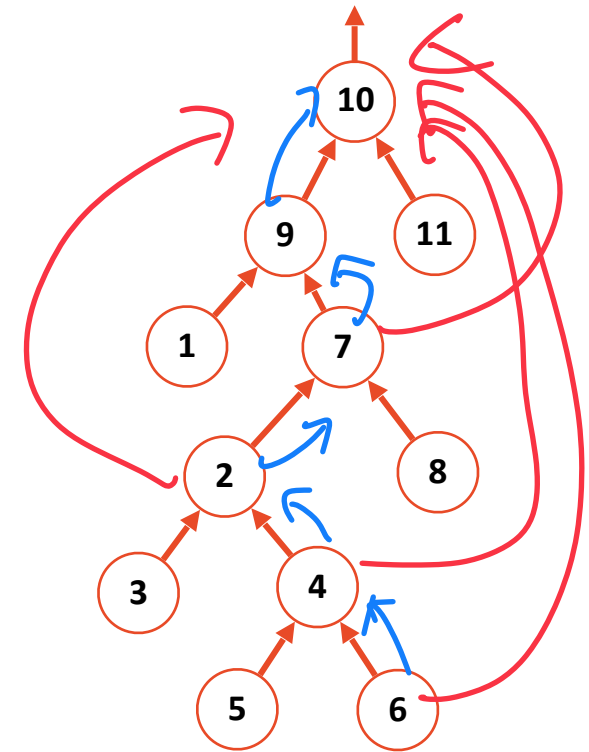
Find(6)

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?

After doing find, update up tree



Claim: why not reach for $\mathcal{O}(1)$?

Disjoint Sets Find

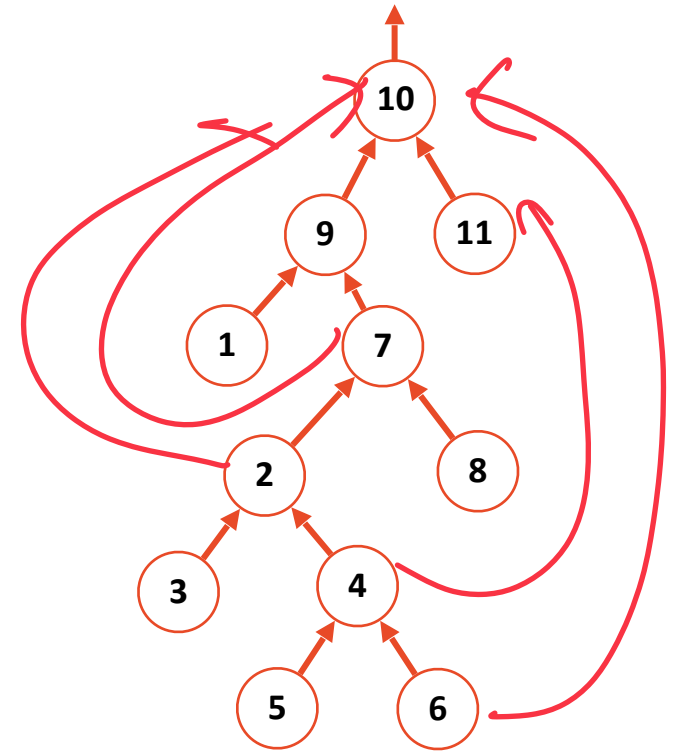
Find(6)

```
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?

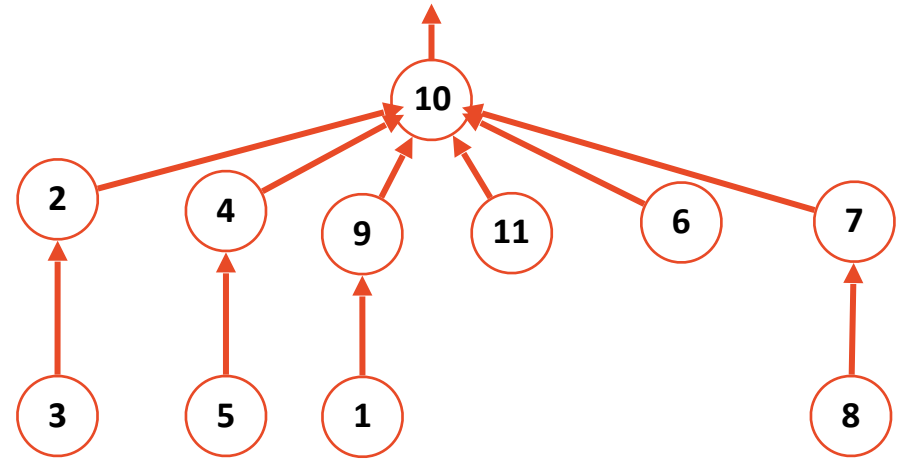
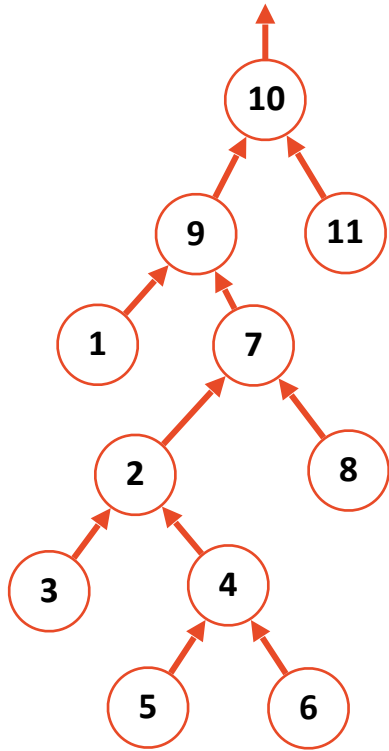
This is **path compression**:

```
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?

Post-Class Edit

We didn't have time to go over this proof in detail!

You should understand path compression but this proof (and rank) is outside scope!

Path Compression Analysis

Two major problems here:

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

↳ Amortized!

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**.

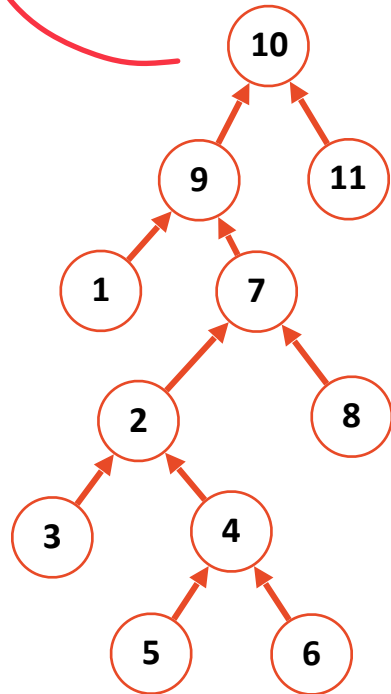


Amortized Time (Path Compression)

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

A.C.

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)

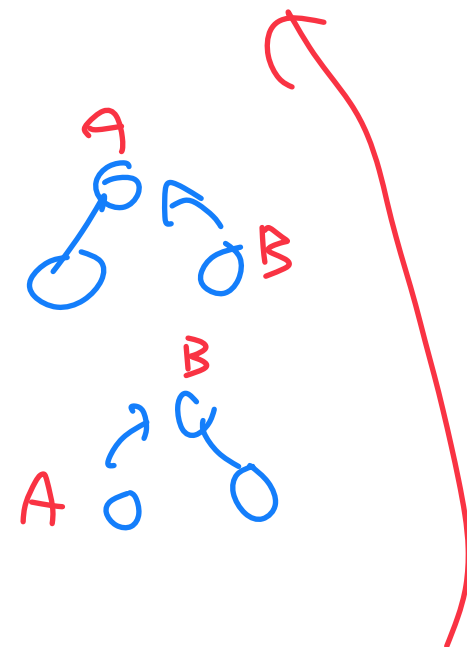
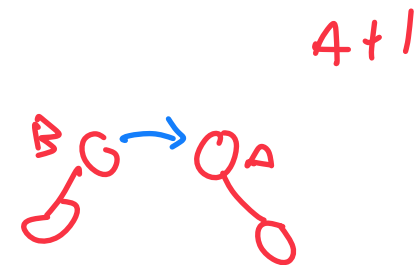
New UpTrees have rank = 0

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)



↳ Rank only increases if we merge 2 sets of same rank

Union by Rank (Not Height)

New UpTrees have rank = 0

A way of putting X nodes under set of rank r,



$r=0$



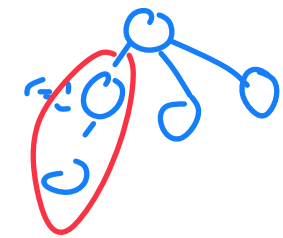
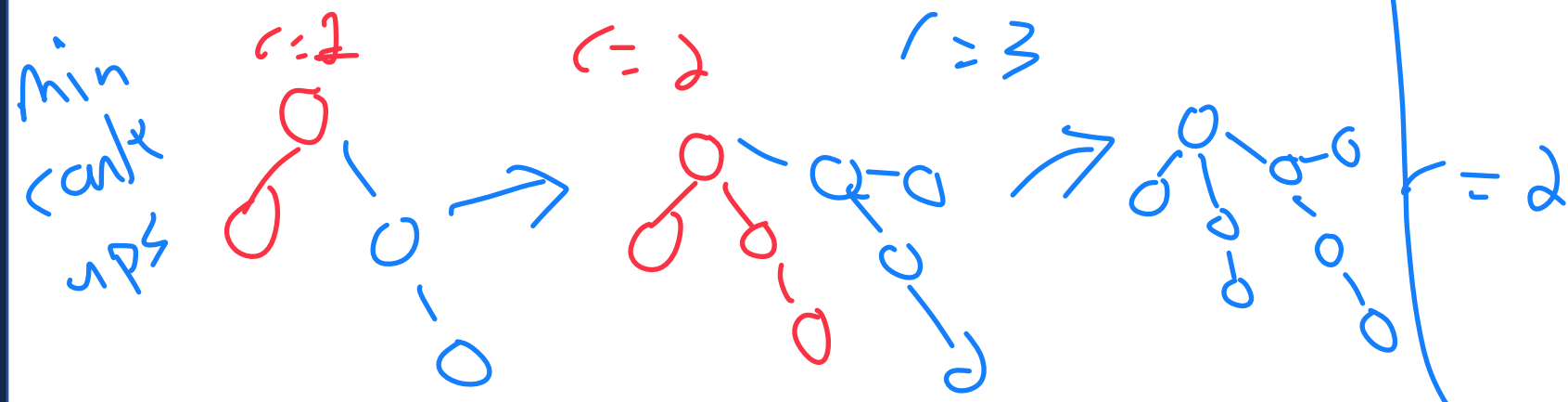
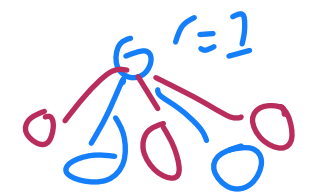
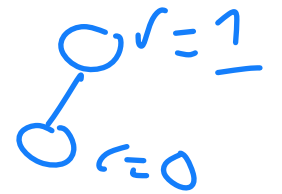
$r=0$

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 $\geq 2^r$ 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)

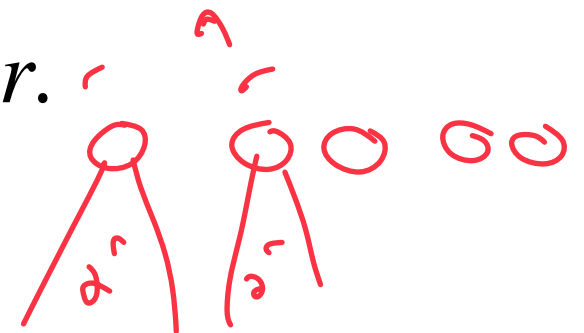
There are at least $\geq 2^r$ nodes in a root of rank r .

rank only increases when union of same rank sets

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

Amortized Time (Rank w/ Path Compression)

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

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

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

Didn't have time for proof - sorry!

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

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

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

$$\log^*(2^{65536}) = 5 = 6$$

Not a true constant but...

CS considers a constant

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

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.

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.

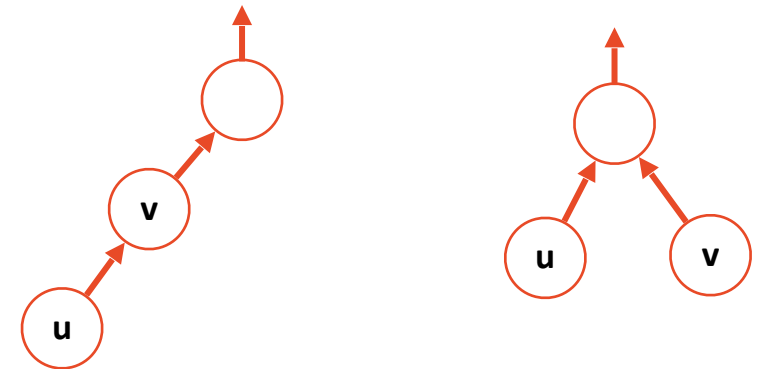
We have at most $\log^*(n)$ buckets so for m finds, this is $O(m \log^* n)$

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



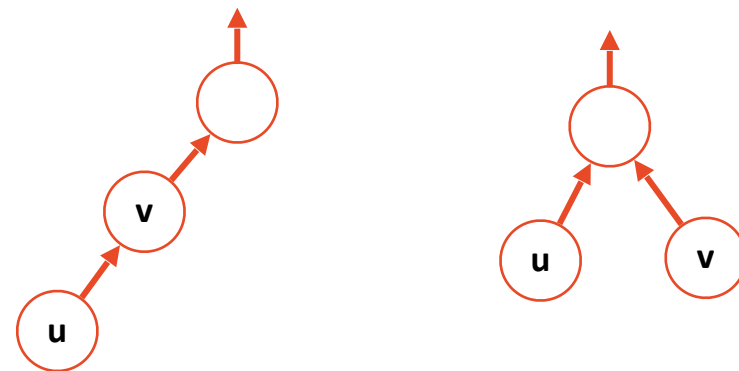
Amortized Time (Rank w/ Path Compression)

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 $\sim 2^r$

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

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

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.

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.