

Union Find

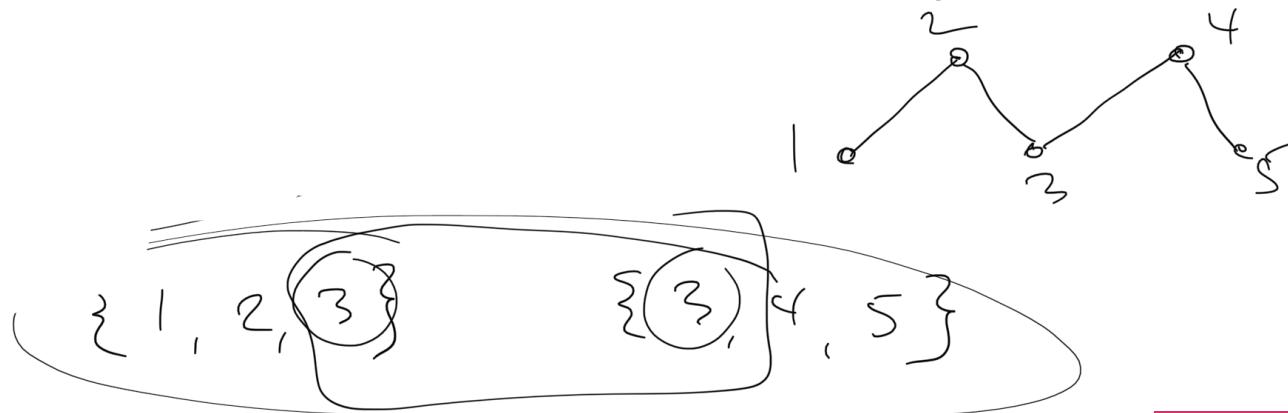
Problem: Dynamic Connectivity

Set - unique items

- Input: a sequence of pairs of integers
- A pair (p, q) means p is connected to q , where connectivity is:
 - reflexive: p is connected to P
 - symmetric: if p is connected to q , then q is connected to p
 - transitive: if p is connected to q and q is conn.-to R , then p is conn. $\xrightarrow{\text{to } R}$
- So connectivity is an equivalence relation, which can separate objects into equivalence classes (i.e. here, two objects are in the same equivalence class if and only if they are connected.)
- Goal: Filter out extraneous pairs (i.e. ignore any pair (p, q) where p and q are ALREADY in the same equivalence class.)

Applications

- mathematical sets (unioning two sets)
- networks → people (graphs)
-



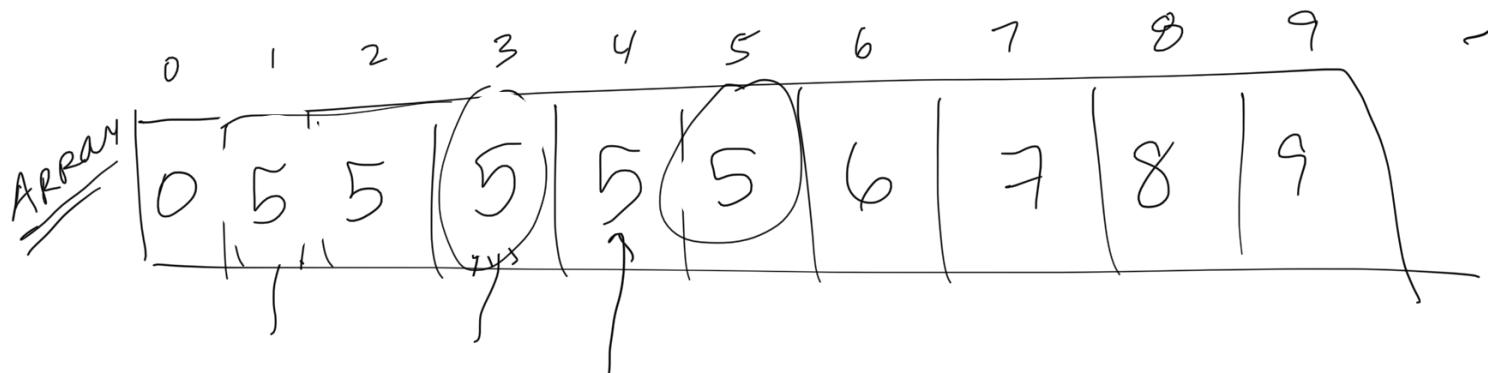
Applications

- Networks
- Variable name equivalence
- Mathematical sets
- Graph connectivity – mST

union (4 3)
union (1 3)
union (1 4)
union (2 5)
union 6

3 is

Quick-Find



Union-Find API

```
public class UF
```

UF(int N)	<i>initialize N sites with integer names (0 to N-1)</i>
void <u>union(int p, int q)</u>	<i>add connection between p and q</i>
int <u>find(int p)</u>	<i>component identifier for p (0 to N-1)</i>
boolean <u>connected(int p, int q)</u>	<i>return true if p and q are in the same component</i>
int <u>count()</u>	<i>number of components</i>

Union-find API

How do we implement this so the operations are efficient?

How do we implement this so the operations are efficient?

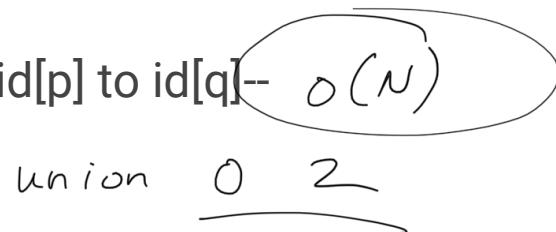
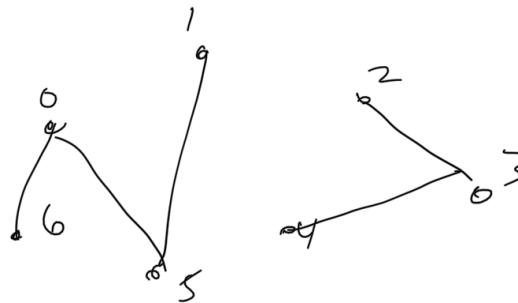
Some general implementation choices...

- The Union-Find implementation will maintain an array called **id** that keeps track of the component of each vertex. (i.e. if $\text{id}[i]=4$, then vertex i is in the component labeled 4.)
- Initially, each vertex is its own component, so $\text{id}[i] = i$ for all i .
- Maintain a count of the number of components. To start with, this is **n** (the number of vertices.)

Option 1: Quick-find

0	1	2	3	4	5	6
0	0	0	0	0	0	0

- UF(int n): initialize id values-- $O(N)$
- union(int p, int q):
 - get $\text{id}[p]$ and $\text{id}[q] \rightarrow O(1)$
 - if ($\text{id}[p] == \text{id}[q]$) do nothing $\rightarrow O(1)$
 - else
 - scan id and set any element that is equal to $\text{id}[p]$ to $\text{id}[q] \rightarrow O(N)$
 - decrement component count
- find(int p): return $\text{id}[p]$ -- $O(1)$
- connected(int p, int q): return $\text{id}[p] == \text{id}[q]$ -- $O(1)$
- count(): return count-- $O(1)$



0

1

2

3

4

5

6

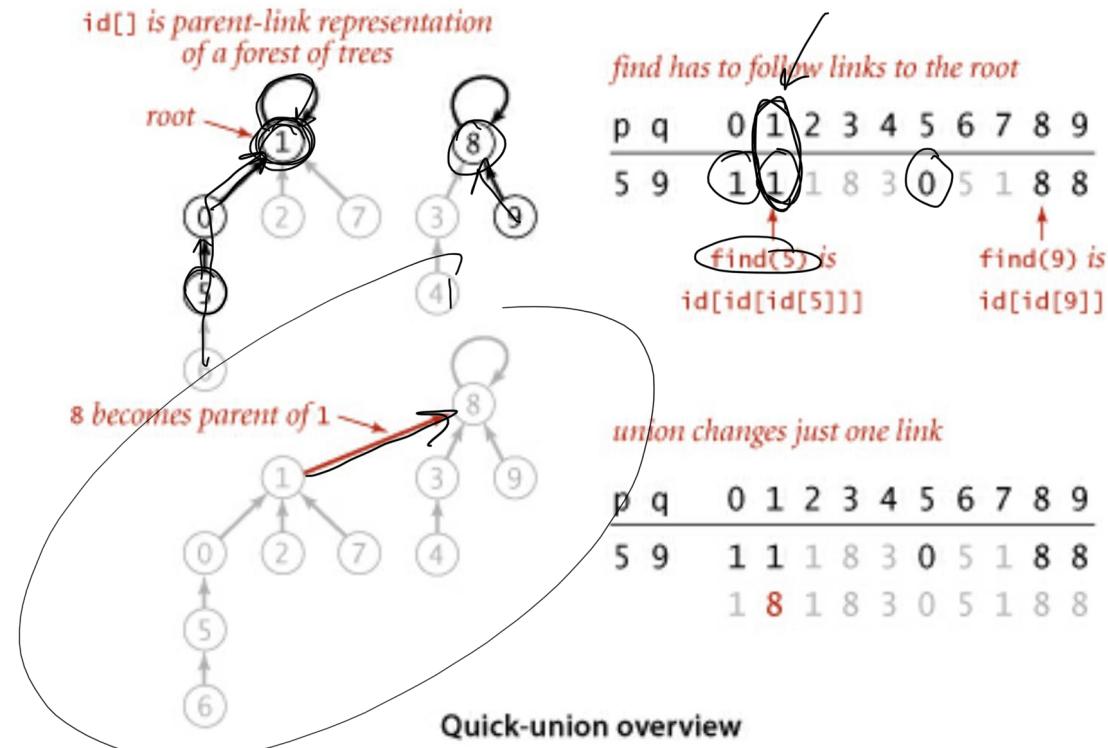
7

8

9

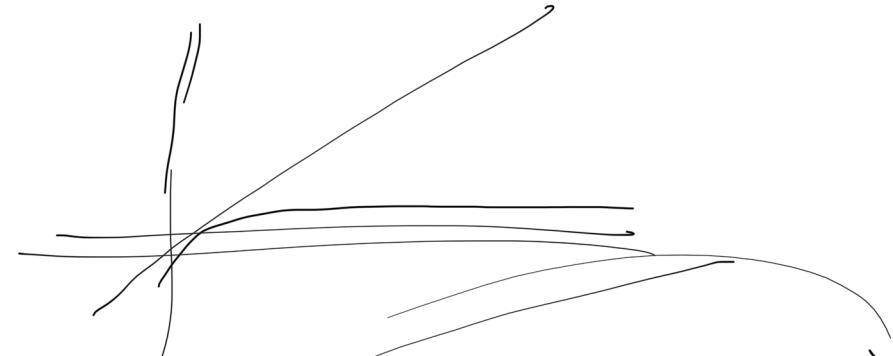
Option 2: Quick-union

- **id** is set up like a tree so that **id[i]** gives you its parent, and so on, until you get to a value that points to itself (the root)
- only one update is needed to union two sets but how many items are checked to find the root?



Option 2: Quick-union

- **UF(int n):** initialize **id** values-- $O(n)$
- **union(int p, int q):**
 - follow the “path” to find the root of the components of **p** and **q**-- $O(\text{tree height})$
 - reassign one root to point to the other
 - decrement component count
- **find(int p):** follow the “path” to find the root-- $O(\text{tree height})$
- **connected(int p, int q):** return **find(p) == find(q)**-- $O(\text{tree height})$
- **count():** return count-- $O(1)$



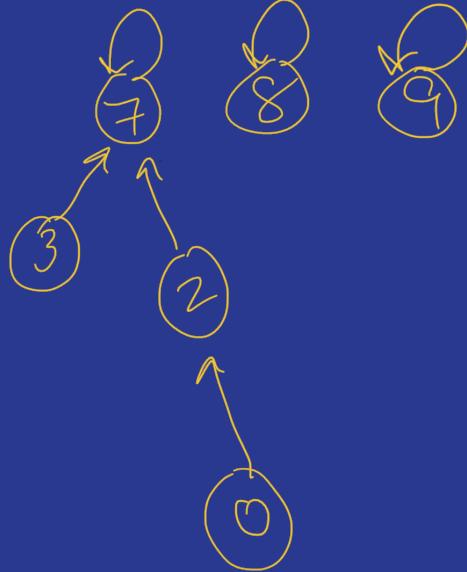
0	1	2	3	4	5	6	7	8	9
2	1	7	7	4	5	6	7	8	9



union 3 7

union 0 2

union 0 7



Count = 7

		id[]									
p	q	0	1	2	3	4	5	6	7	8	9
4	3	0	1	2	3	4	5	6	7	8	9
		0	1	2	3	3	5	6	7	8	9
3	8	0	1	2	3	3	5	6	7	8	9
		0	1	2	8	3	5	6	7	8	9
6	5	0	1	2	8	3	5	6	7	8	9
		0	1	2	8	3	5	5	7	8	9
9	4	0	1	2	8	3	5	5	7	8	9
		0	1	2	8	3	5	5	7	8	8
2	1	0	1	2	8	3	5	5	7	8	8
		0	1	1	8	3	5	5	7	8	8
8	9	0	1	1	8	3	5	5	7	8	8
5	0	0	1	1	8	3	5	5	7	8	8
		0	1	1	8	3	0	5	7	8	8
7	2	0	1	1	8	3	0	5	7	8	8
		0	1	1	8	3	0	5	1	8	8
6	1	0	1	1	8	3	0	5	1	8	8
		1	1	1	8	3	0	5	1	8	8
1	0	1	1	1	8	3	0	5	1	8	8
6	7	1	1	1	8	3	0	5	1	8	8

Quick-union trace (with corresponding forests of trees)

Quick-union trace with representation as a forest of trees.

As connections are added, you get fewer but larger trees (correspond to components).

If the runtime of key operations depends on the height of the tree, what is the worst case?

		id[]					
p	q	0	1	2	3	4	...
0	1	0	1	2	3	4	...
		1	1	2	3	4	...
0	2	0	1	2	3	4	...
		1	2	3	4	...	
0	3	0	1	2	3	4	...
		1	2	3	3	4	...
0	4	0	1	2	3	4	...
		1	2	3	4	4	...
.							
.							
.							

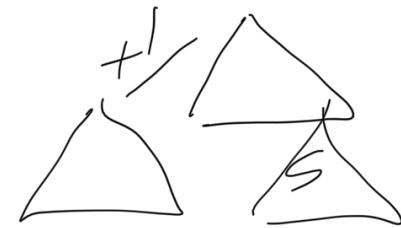
depth 4 → 0

Quick-union worst case

Worst Case: O(N)

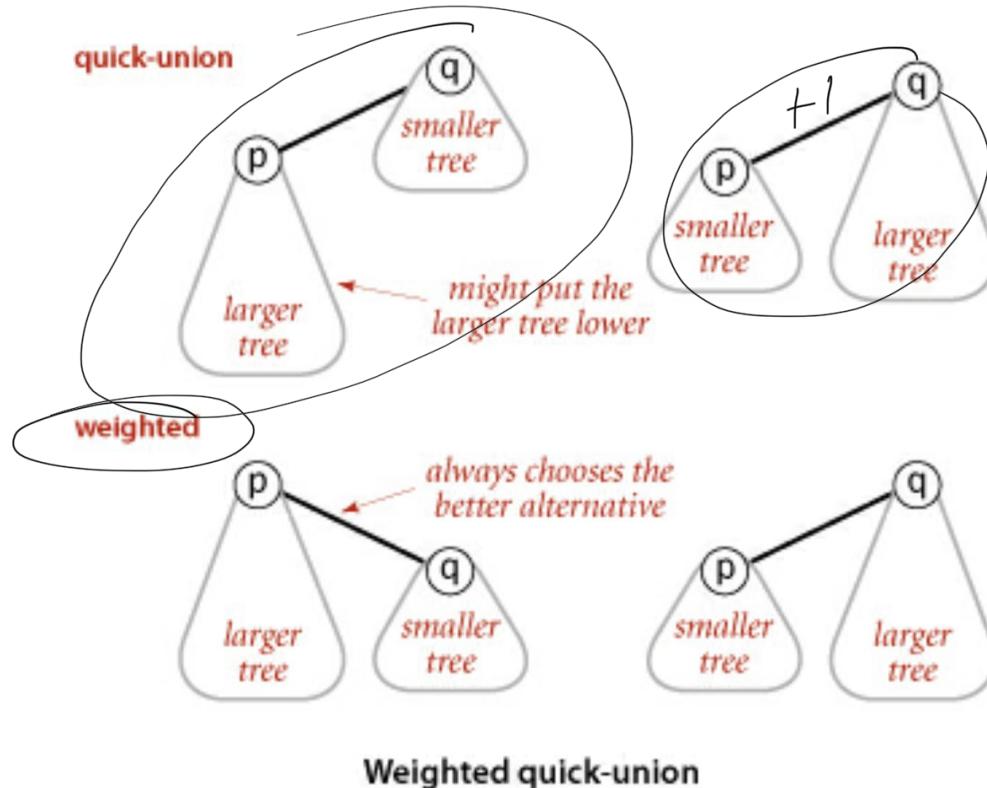
Can we improve upon this?

Option 3: Weighted Quick-union



Analysis:

- In the worst case, the trees are of equal size and combining them results in a tree of twice the size whose height increases by 1.
- Consider two trees with 2^k nodes. Then combining those trees will produce a new tree of size 2^{k+1} where the height has increased by 1.
- If we produce all trees this way, we guarantee that the height of the tree is $O(\log N)$.



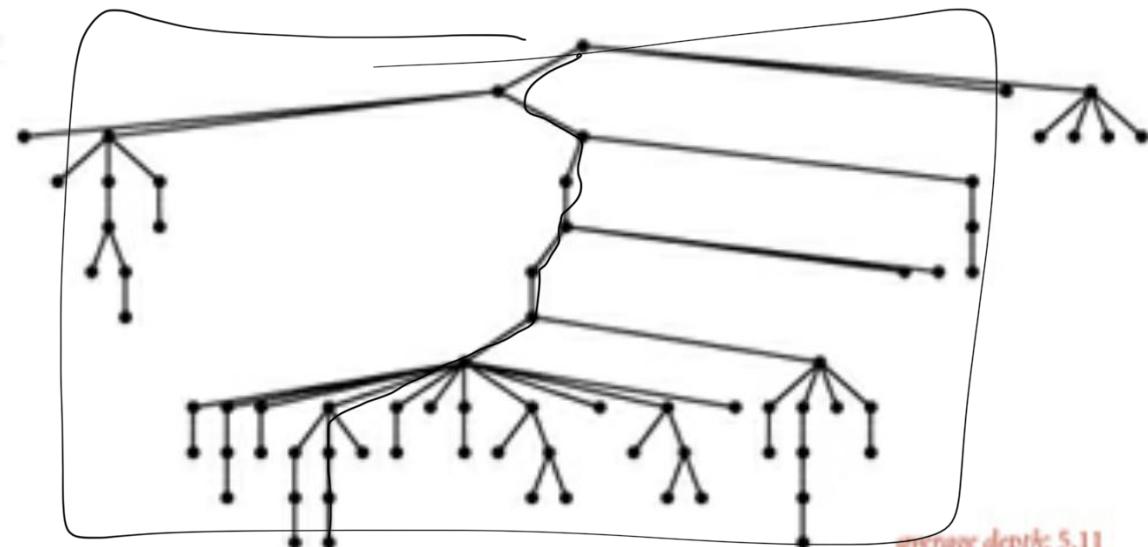
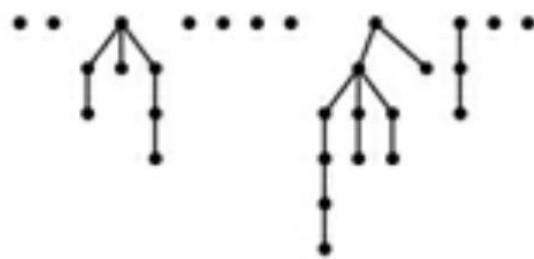
Weighted quick-union

Option 3: Weighted Quick-union

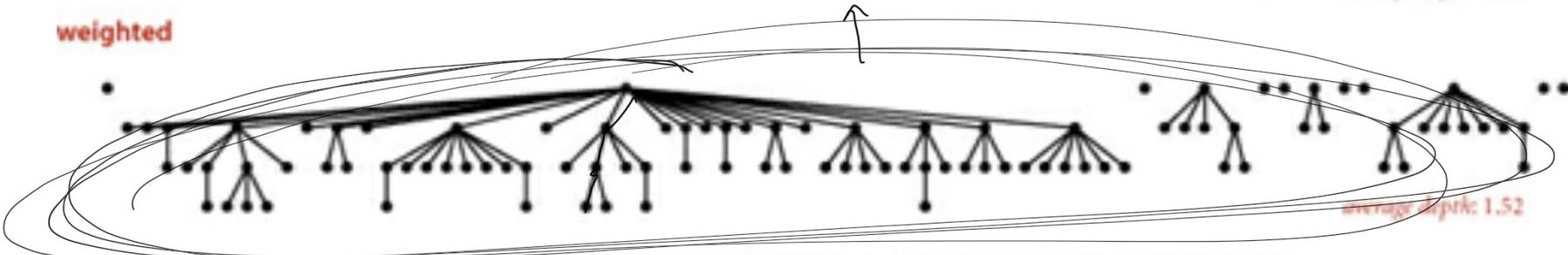
- **UF(int n):** initialize **id** values-- $O(N)$
- **union(int p, int q):**
 - follow the “path” to find the root of the components of **p** and **q**-- $O(\log N)$
 - reassign the root of the smaller tree to point to the other
 - decrement component count
- **find(int p):** follow the “path” to find the root-- $O(\log N)$
- **connected(int p, int q):** return **find(p) == find(q)**-- $O(\log N)$
- **count():** return count-- $O(1)$



quick-union



weighted



Quick-union and weighted quick-union (100 sites, 88 union() operations)

0

1

2

3

4

5

6

7

8

9

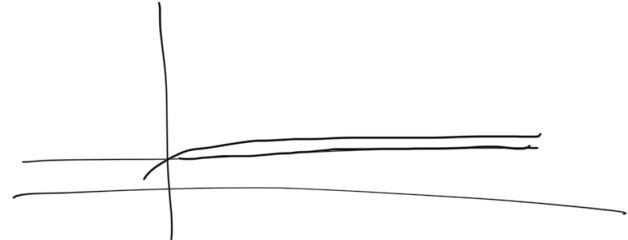
Option 4: Weighted Quick-union with Path Compression



- Weighted Quick-Union with an extra optimization.
- **Main Idea:** Ideally, we want every node to link directly to its root.
- **But...** it's expensive to change all the nodes (remember Quick-Find?)
- **Solution:** Change the ones you are already examining as you look for the root--this just requires another loop to the *find* operation that sets the `id[]` value of each node on the path to the root. What's the runtime?

Option 4: Weighted Quick-union with Path Compression

What's the runtime?



- The worst-case for *find* and *union* may still be $O(\log N)$ but theoretical analysis shows that the amortized time for both is **close to constant**.
- However, it should be noted that you're unlikely to see much of a difference in practical situations between this and simple weighted quick-union.

algorithm	<i>order of growth for N sites (worst case)</i>		
	constructor	union	find
quick-find	N	N	1
quick-union	N	tree height	tree height
weighted quick-union	N	$\lg N$	$\lg N$
weighted quick-union with path compression	N	<u>very, very nearly, but not quite 1 (amortized)</u> (see EXERCISE 1.5.13)	
impossible	N	1	1

Performance characteristics of union-find algorithms

References

[1] *Algorithms, Fourth Edition*; Robert Sedgewick and Kevin Wayne