

CSCI 335

Software Design and Analysis III

Lecture 20:

The Disjoint Sets Class

Professor Anita Raja

11-14-22

Agenda

- Equivalence Relations
 - Union-Find
 - Smart-Union
 - Union by size
 - Union by height
 - Path compression
 - Union by rank
- Analysis
- Application

Equivalence Relations

- **Relation R** on a set S: for every pair (a,b) , $a, b \in S$, $a R b$ is either T or F.
 - If $a R b$ is T, we say a is related to b.
 - Examples:
 - Relation $<$ on the set $\{0,1,2,3,4,5,6\}$
 -Many more....
- **Equivalence** relation R has 3 properties
 - $a R a$, for all a (reflexive)
 - $a R b \Rightarrow b R a$ (symmetric)
 - $a R b$ and $b R c \Rightarrow a R c$ (transitive)
- **Examples** (which ones of the following are equivalence relations?):
 - $=$?
 - \leq ?
 - Two cities in the same country ?

Dynamic Equivalence Problem

- Given set $S=\{a_1,a_2,a_3,a_4,a_5,a_6\}$
 - 36 pairs of elements that are either related or not
- Given equivalence relation \sim
 - $a_1 \sim a_2, a_2 \sim a_3, a_5 \sim a_6$
 - $a_i \sim a_i$ for $i=1,\dots,6$ (from reflexivity)
 - Symmetry $\Rightarrow a_2 \sim a_1, a_3 \sim a_2, a_6 \sim a_5$
 - Transitivity $\Rightarrow a_1 \sim a_3, a_3 \sim a_1$
- **Equivalence class** of an element $a \in S$: is subset of $S \in$ all elements related to a
 - The equivalence classes form a partition of S .
 - Every member of S appears in exactly one equivalence class.
 - To decide if $a \sim b$, we need only check whether a and b are in the same equivalence class
 - $\{a_1,a_2,a_3\}, \{a_4\}, \{a_5,a_6\}$ are the equivalence classes in above example

Disjoint sets

- Input is initially a collection of N sets, each with one element.
- Initial representation is that all relations (except reflexive relations) are false.
- Each set has a different element, so that $S_i \cap S_j = \emptyset$

=> Disjoint sets

Union-Find

- Find:

- given an element, find the equivalence class it is in

- Union:

- Add the relation $a \sim b$:
 - Perform a Find on a and b.
 - If already related, do nothing.
 - Else Union the equivalence classes of a and b into a new class.
- $S_i \cup S_j = S_k$

Union-Find

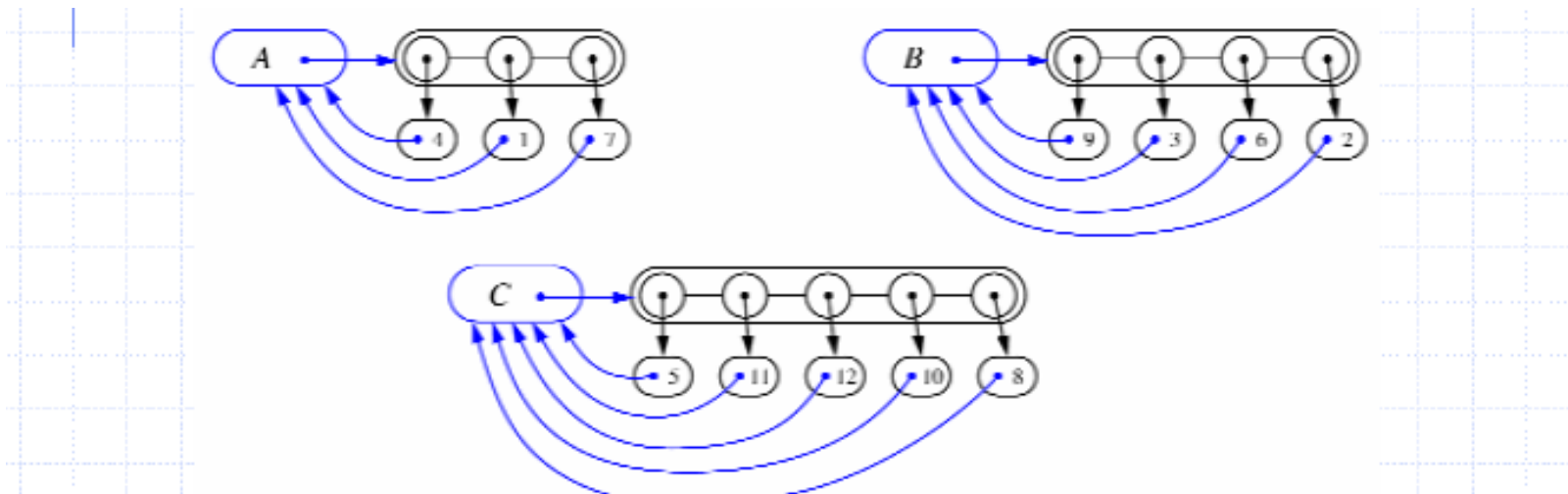
- **Dynamic:** the equivalence classes can change with time,
 - i.e. Unions are applied at any time during the course of the algorithm.
- **On-line:** a union or find “arrives” at each instance of time and needs to be processed.
 - i.e. when a find is performed, it must give an answer before continuing.
 - Like an oral exam
- **Off-line:** the sequence of unions and finds is given, and you need to calculate final sets (not dynamic).
 - Like a written exam

Union-Find

- Assume set $S = \{0,1,2,\dots,N-1\}$ (can be hashed)
- Name of set returned by find is irrelevant
- Solutions?
- $O(1)$ find?
 - Union ?

List-based Implementation

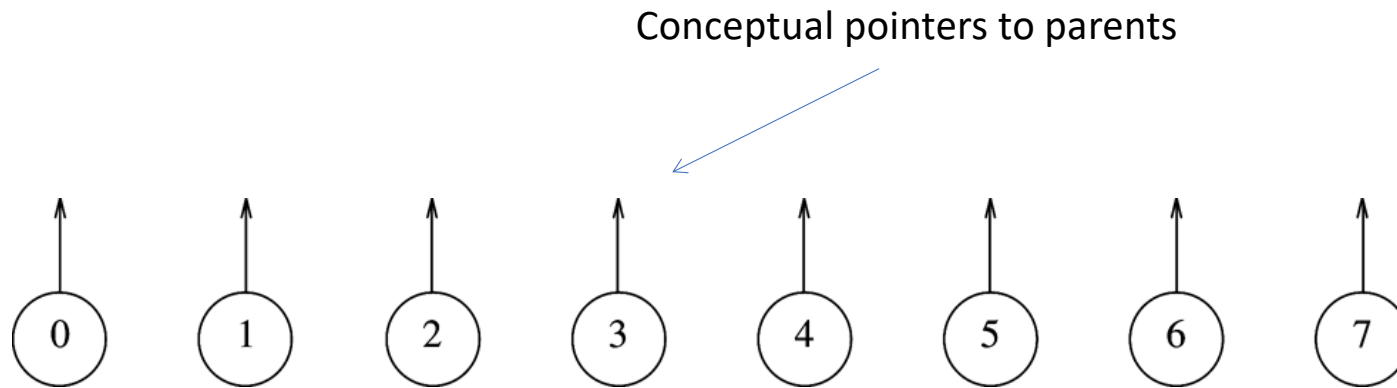
- Each set is stored in a sequence represented with a linked-list.
- Each node should store an object containing the element and a reference to the set name



Analysis of List-based representation

- When doing a union, always move elements from the smaller set to the larger set
 - Each time an element is moved it goes to a set of size at least double its old set.
 - Thus, an element can be moved at most $O(\log N)$ times.
- Total time needed to do M finds and $N-1$ unions is $O(M+N\log N)$
- $O(N\log N)$

Better



Each equivalence class is represented by a tree.

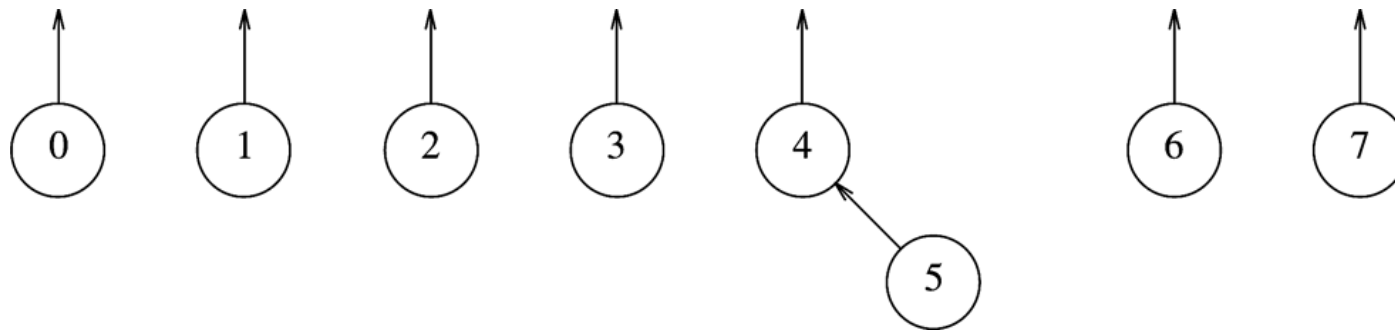
The name of the class is the root (i.e. name of the representative).

Initially each element is one class containing itself (reflexivity).

Can be implemented by array

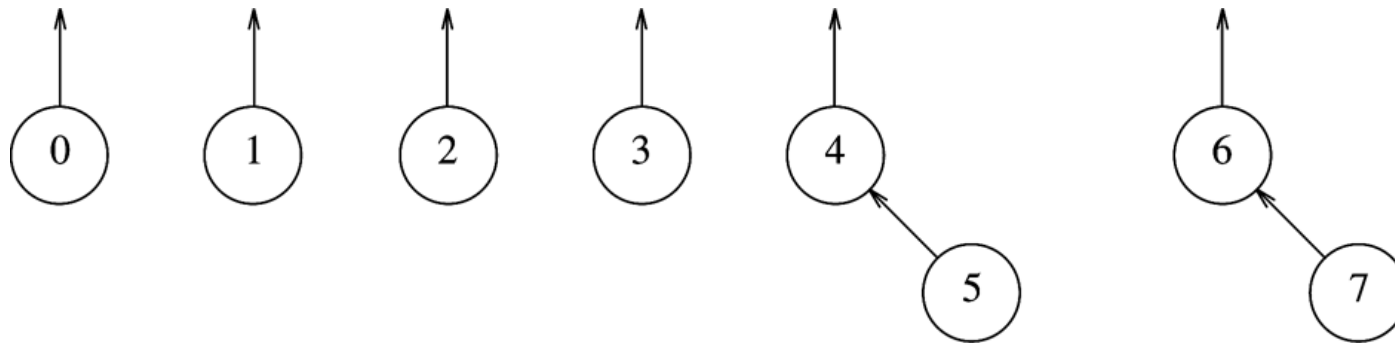
Better

After Union(4,5)



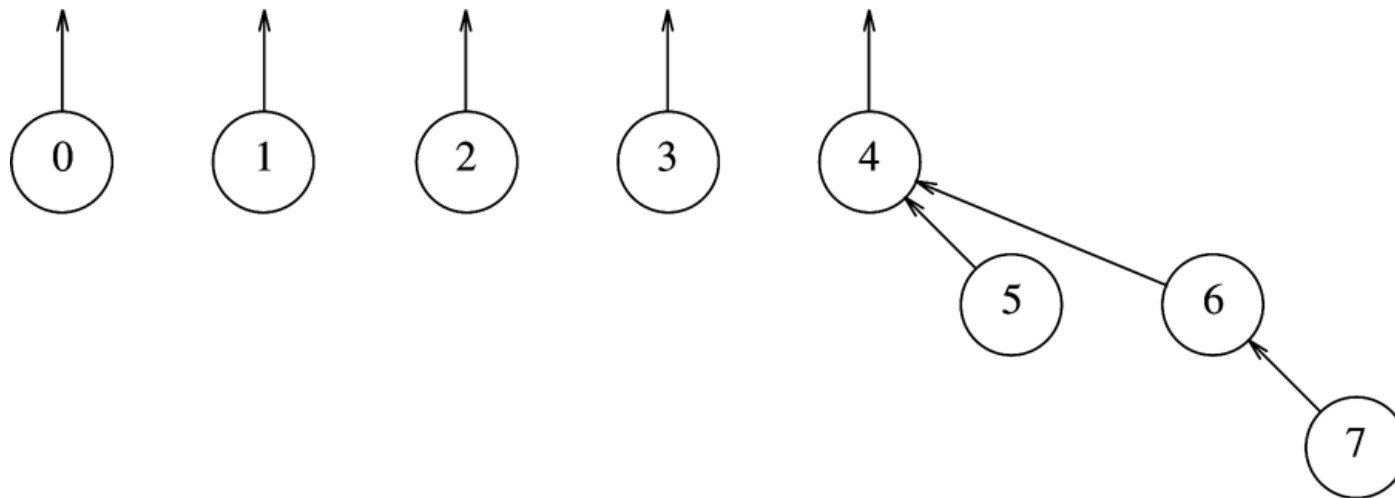
Better

After Union(6,7)



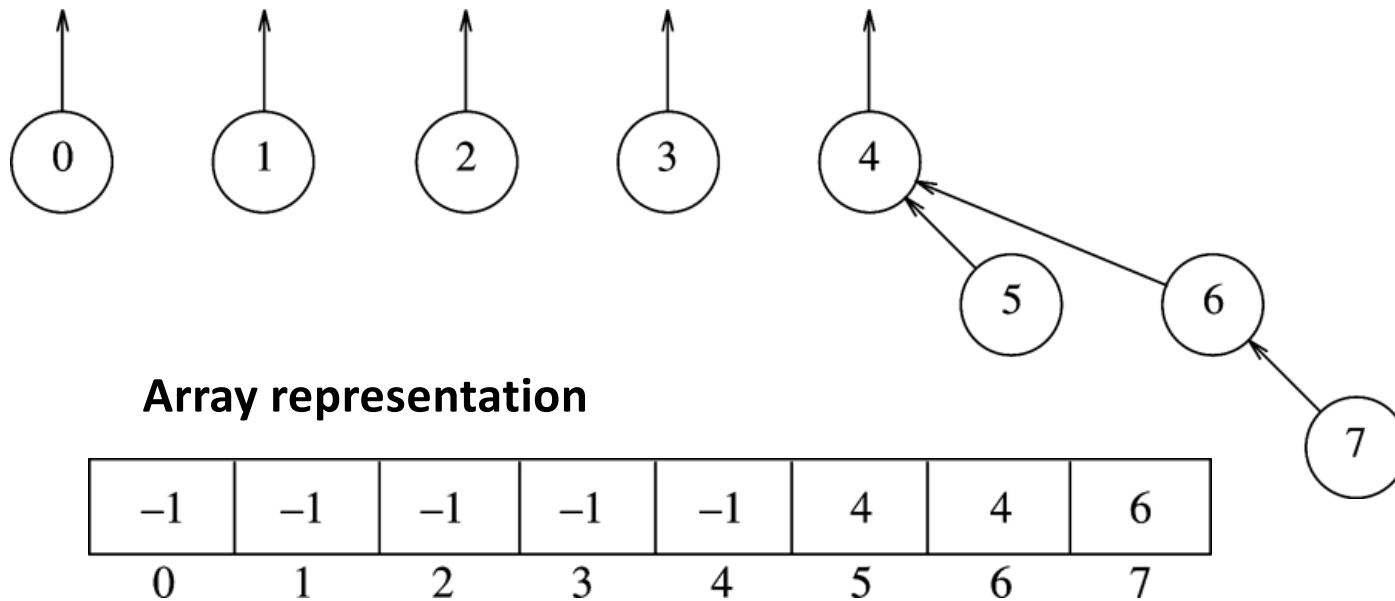
Better

After Union(4,6)



Better

After Union(4,6)



Implementation

```
1  class DisjSets
2  {
3      public:
4          explicit DisjSets( int numElements );
5
6          int find( int x ) const;
7          int find( int x );
8          void unionSets( int root1, int root2 );
9
10     private:
11         vector<int> s;
12 };
```


Analysis

- Union's cost?

Analysis

- Union's cost?
- Union: Fast $O(1)$

Analysis

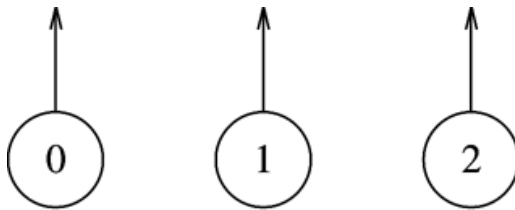
- Union's cost?
- Union: Fast $O(1)$
- Find's cost?

Analysis

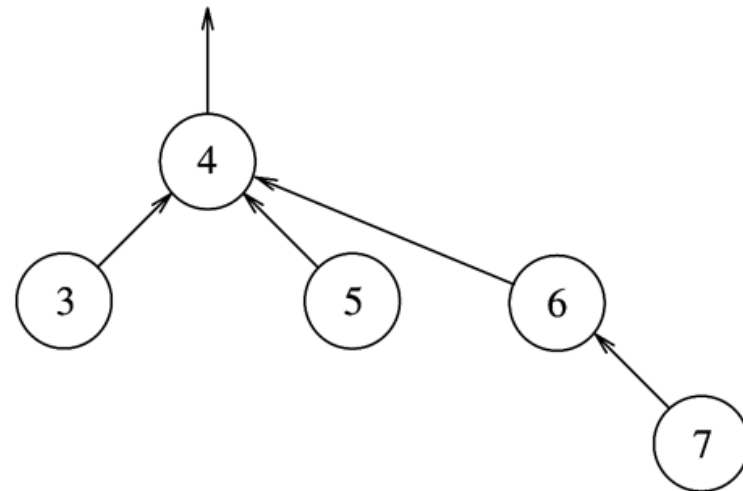
- Union's cost?
- Union: Fast $O(1)$
- Find's cost ?
- Find: Worst-case depth of deeper node $O(N)$
=> $O(MN)$ for M mixed find/union operations
- Average running time: depends on what is *average*:
 - $\Theta(M)$ or $\Theta(MN)$ or $\Theta(M \log N)$ for M unions

Smart Union Algorithms

- Union by size: Smaller tree becomes subtree of larger



After Union(3,4)



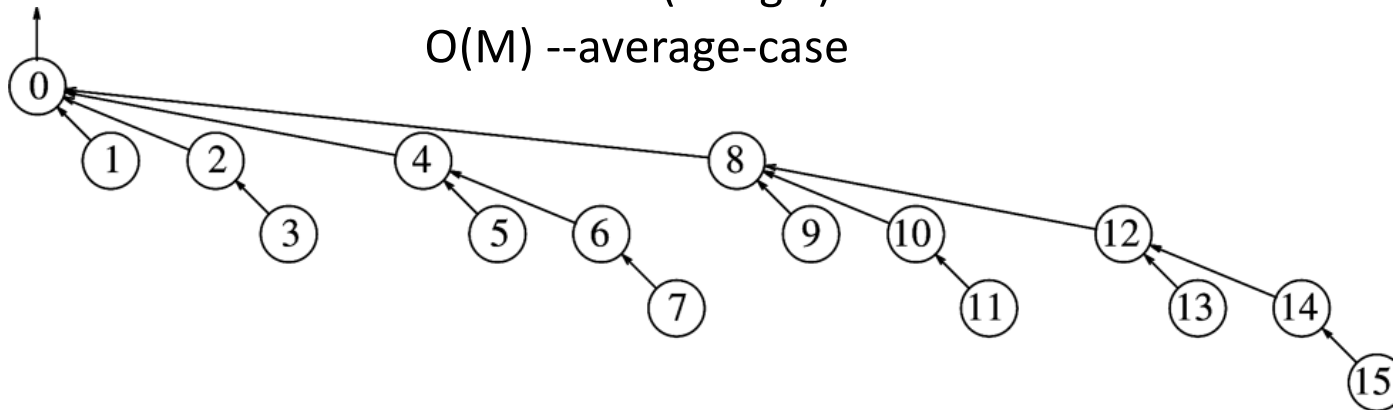
Smart Union Algorithms

Depth of any node can't be more than $\log N$

\Rightarrow Find costs $O(\log N)$ --worst-case

\Rightarrow M finds cost $O(M \log N)$ --worst-case

$O(M)$ --average-case



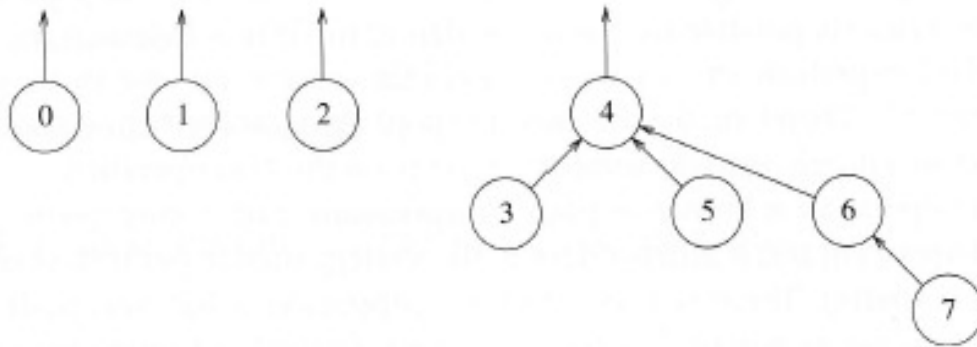
Worst-case for union when $N=16$ (binomial tree) [result after 16 unions]

[union of equal size trees]

Smart Union Algorithms

- **Union-by-height:** shortest tree under tallest
- Height is increased by one **iff** when two trees of equal height are combined.
- Again $O(\log N)$ find

Smart Union Algorithms



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

Union-by-size

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

Union-by-height

Smart Union Algorithms

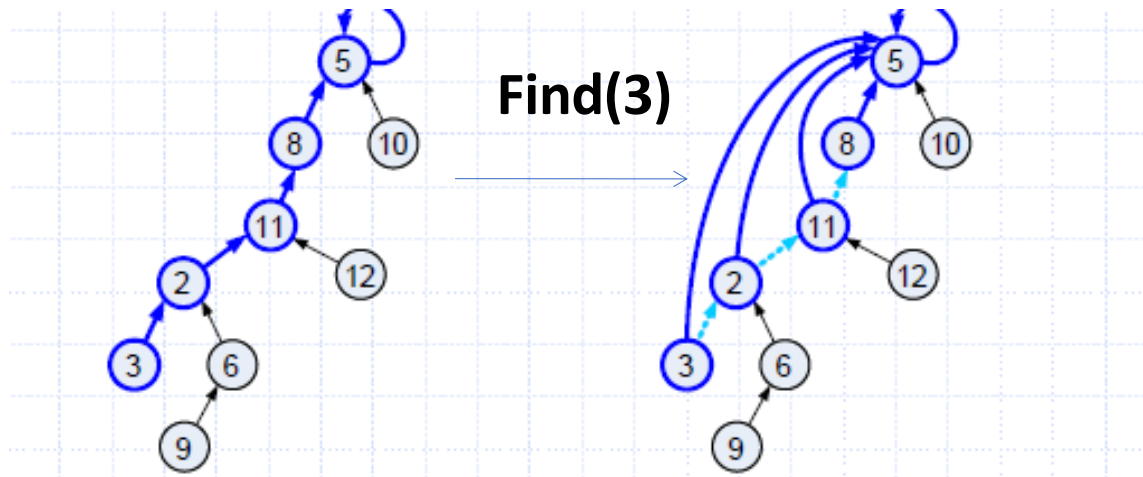
```
1  /**
2   * Union two disjoint sets.
3   * For simplicity, we assume root1 and root2 are distinct
4   * and represent set names.
5   * root1 is the root of set 1.
6   * root2 is the root of set 2.
7   */
8  void DisjSets::unionSets( int root1, int root2 )
9  {
10     if( s[ root2 ] < s[ root1 ] ) // root2 is deeper
11         s[ root1 ] = root2;      // Make root2 new root
12     else
13     {
14         if( s[ root1 ] == s[ root2 ] )
15             s[ root1 ]--;        // Update height if same
16         s[ root2 ] = root1;      // Make root1 new root
17     }
18 }
```

Path Compression

- Improve on the $O(M \log N)$ worst-case `find()` cost for M finds
- Ideas?

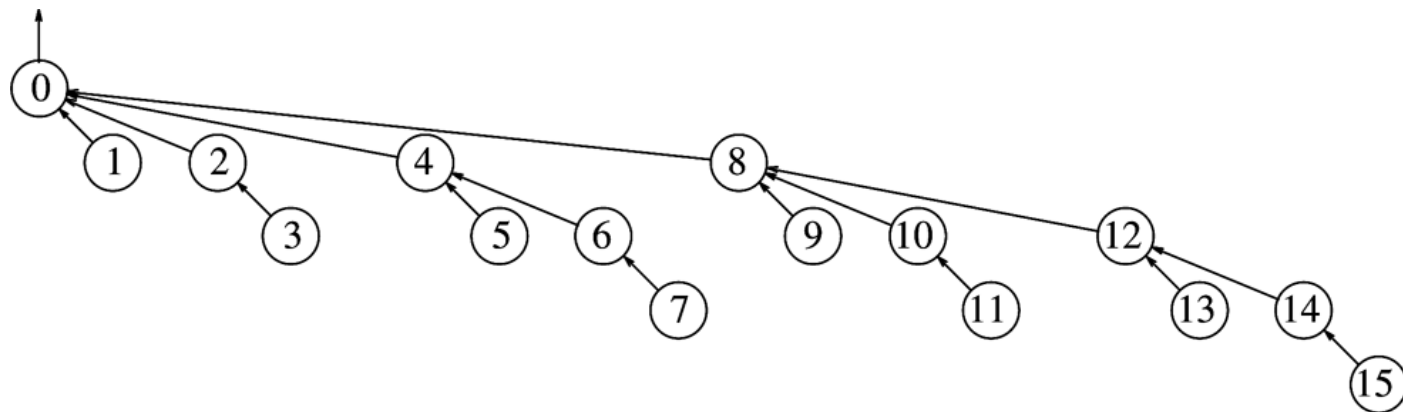
Path Compression

- Improve on the $O(M \log N)$ worst-case $\text{find}()$ cost for M finds
- It can appear frequently
- **Path compression:** During $\text{find}(x)$ make all nodes on the path from x to the root to be children of the root.

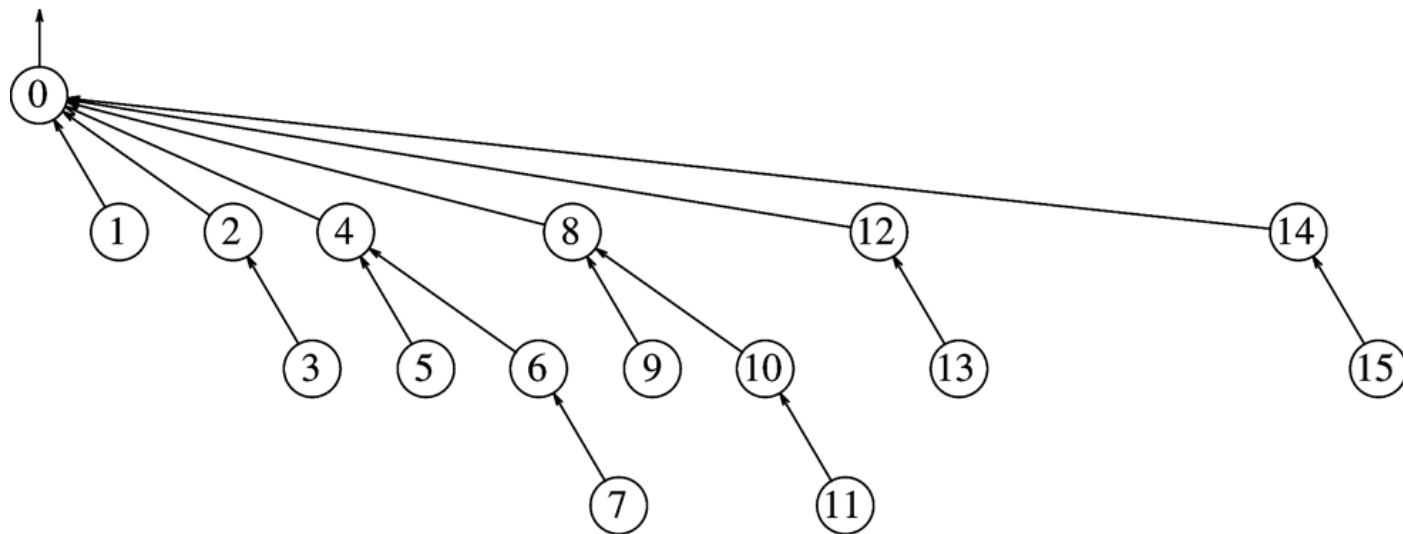


Path Compression

Apply find (14) in this set



Path Compression



Path Compression

```
1  /**
2   * Perform a find with path compression.
3   * Error checks omitted again for simplicity.
4   * Return the set containing x.
5   */
6  int DisjSets::find( int x )
7  {
8      if( s[ x ] < 0 )
9          return x;
10     else
11         return s[ x ] = find( s[ x ] );
12 }
```

Union-by-rank

- **Union-by-rank:** after compression the height of tree may not be accurate.
- It is now called **rank** => it is an upper bound on the actual height.

Path Compression

(A) Union-by-size + Path Compression

or

(B) Union-by-rank + Path Compression

Average case: Not sure whether Path Compression helps

Worst case: Path Compression helps a lot

Union-by-rank is preferred because it requires fewer updates on heights

Worst-Case Analysis of Union-Find

- Definition: \log^*N is number of times logarithm of N needs to be applied until it gets to ≤ 1 .
- For example: $\log^*65536=4$, because $(65536=2^{16})$
 $\log \log \log \log 65536 = 1$
- Example: $\log^*2^{65536} = ?$ (2^{65536} has 20,000 digits in decimal form)
- \log^*N grows very slowly as N becomes bigger !
 - $\log^*2 = 1$
 - $\log^*4 = 2$
 - $\log^*16 = 3$
 - $\log^*65536 = 4$
 - $\log^*2^{65536} = 5$

Analysis of Union-Find

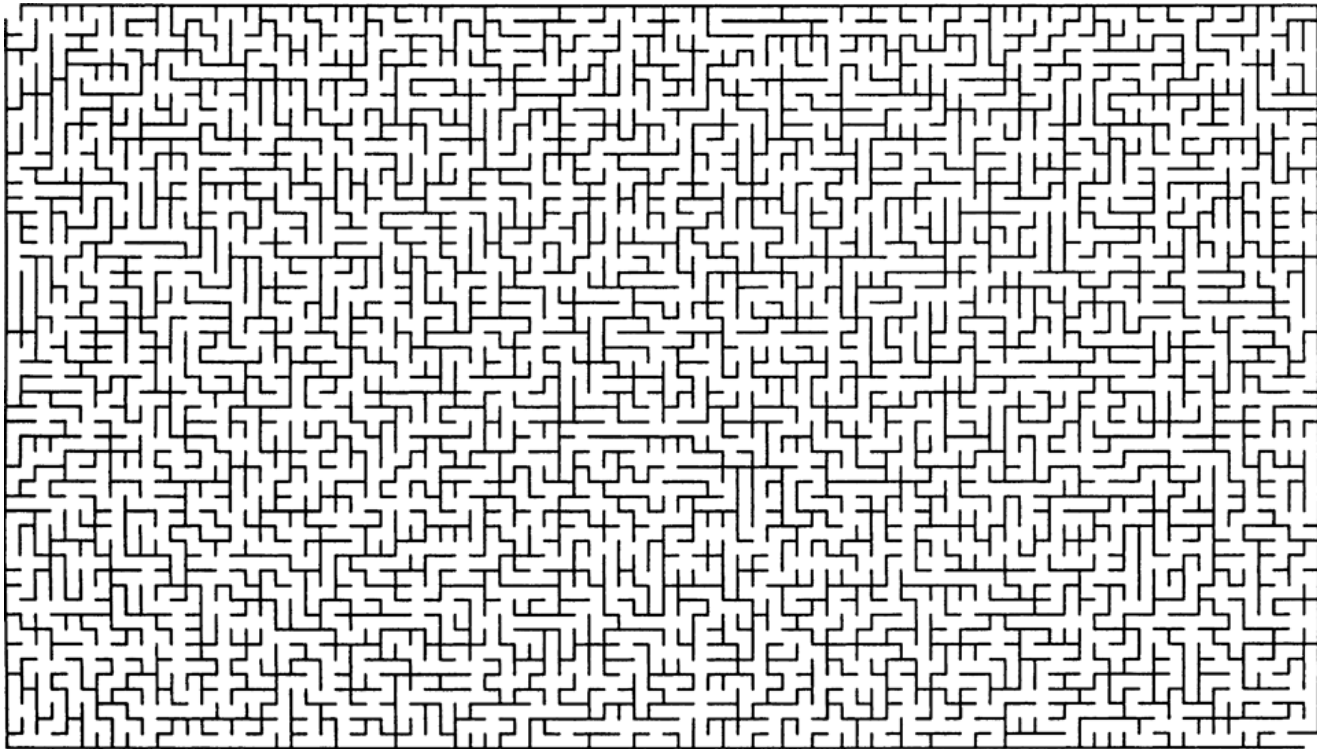
Any sequence of $M = \Omega(N)$ union/find operations takes a total of $O(M \log^* N)$ running time.

Model:

Union/Finds in any order

Union-by-rank **with** path compression

A fun application



A fun application

0	1	2	3	4
5	6	7	8	9
10	11	12	13	14
15	16	17	18	19
20	21	22	23	24

{0} {1} {2} {3} {4} {5} {6} {7} {8} {9} {10} {11} {12} {13} {14} {15} {16} {17} {18} {19} {20} {21}
{22} {23} {24}

A fun application

0	1	2	3	4
5	6	7	8	9
10	11	12	13	14
15	16	17	18	19
20	21	22	23	24

{0, 1} {2} {3} {4, 6, 7, 8, 9, 13, 14} {5} {10, 11, 15} {12} {16, 17, 18, 22} {19} {20} {21} {23} {24}

A fun application

0	1	2	3	4
5	6	7	8	9
10	11	12	13	14
15	16	17	18	19
20	21	22	23	24

{0, 1} {2} {3} {4, 6, 7, 8, 9, 13, 14, 16, 17, 18, 22} {5} {10, 11, 15} {12} {19} {20} {21} {23} {24}

A fun application

0	1	2	3	4
5	6	7	8	9
10	11	12	13	14
15	16	17	18	19
20	21	22	23	24

{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24}

Summary

- Simple data structure to maintain disjoint sets
- Important for graph theoretical problems
- Union step is flexible and so we get a much more effective algorithm
- Path compressions – earliest forms of self-adjustment
 - Seen in splay trees and skew heaps
 - Use extremely interesting because from theoretical pov one of the first algorithms that was simple with a not-so-simple worst case analysis.
- Any sequence of $M = \Omega(N)$ union/find operations takes a total of $O(M \log^* N)$ running time.

Model: Union/Finds in any order

Union-by-rank **with** path compression