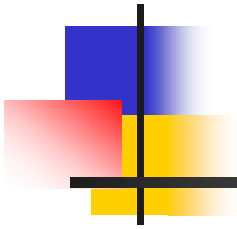


# Union-Find: A Data Structure for Disjoint Set Operations





# The Union-Find Data Structure

- Purpose:
  - To manipulate disjoint sets (i.e., sets that don't overlap)
  - Operations supported:

<i><b>Union ( x, y )</b></i>	Performs a union of the sets containing two elements x and y
<i><b>Find ( x )</b></i>	Returns a pointer to the set containing element x

Q) Under what scenarios would one need these operations?



# Some Motivating Applications for Union-Find Data Structures

---

*Given a set  $S$  of  $n$  elements,  $[a_1 \dots a_n]$ , compute all its equivalent classes*

Example applications:

- Electrical cable/internet connectivity network
- Cities connected by roads
- Cities belonging to the same country



# Equivalence Relations

---

- An equivalence relation  $R$  is defined on a set  $S$ , if for every pair of elements  $(a,b)$  in  $S$ ,
  - $a R b$  is either false or true
- $a R b$  is true iff:
  - (Reflexive)  $a R a$ , for each element  $a$  in  $S$
  - (Symmetric)  $a R b$  if and only if  $b R a$
  - (Transitive)  $a R b$  and  $b R c$  implies  $a R c$
- The equivalence class of an element  $a$  (in  $S$ ) is the subset of  $S$  that contains all elements *related to  $a$*

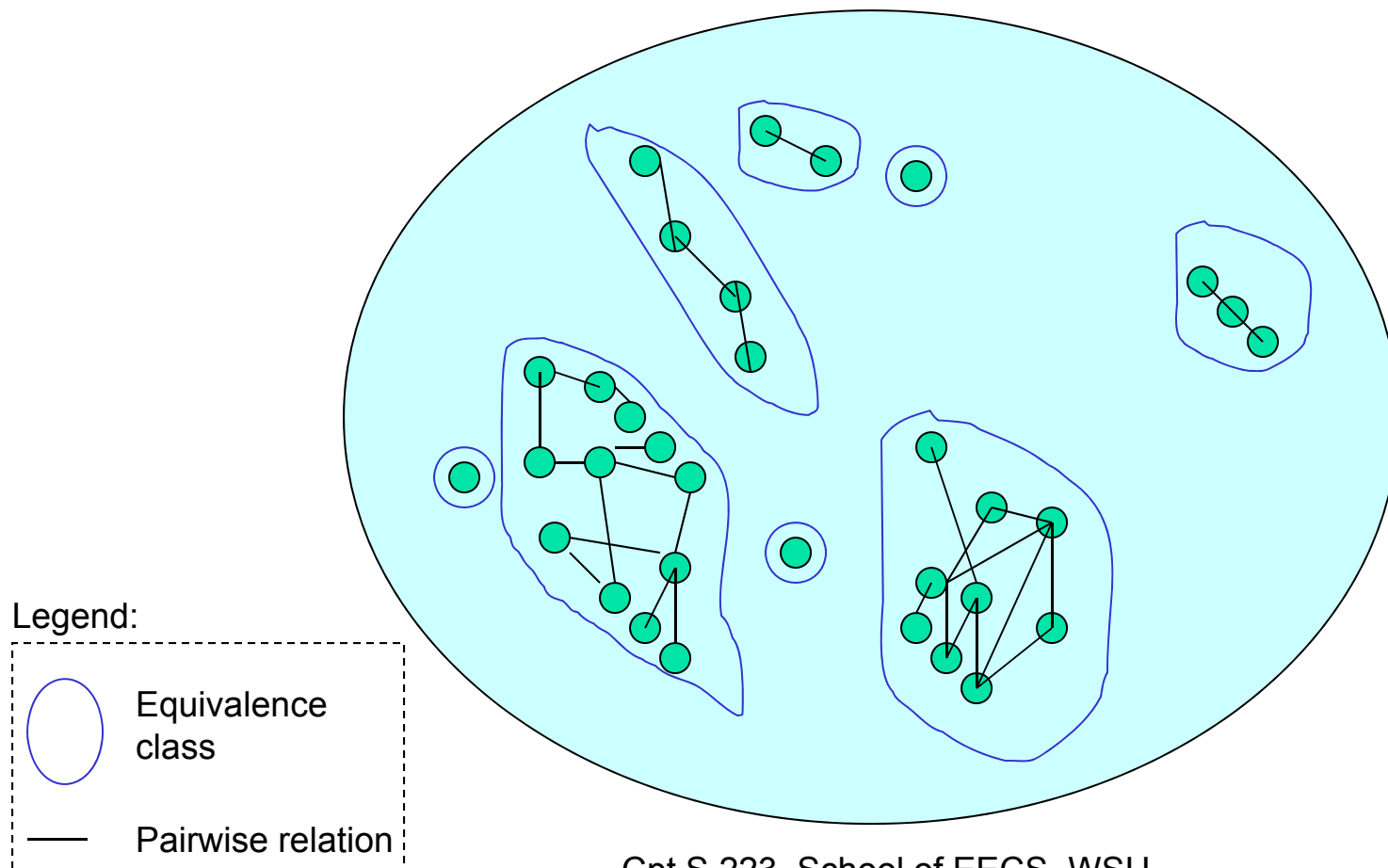


# Properties of Equivalence Classes

---

- An observation:
  - Each element must belong to exactly one equivalence class
- Corollary:
  - All equivalence classes are mutually disjoint
- What we are after is the set of all equivalence classes

# Identifying equivalence classes





# Disjoint Set Operations

---

- To identify all equivalence classes
  1. Initially, put each each element in a set of its own
  2. Permit only two types of operations:
    - **Find(x)**: Returns the current equivalence class of x
    - **Union(x, y)**: Merges the equivalence classes corresponding to elements x and y (assuming x and y are related by the eq.rel.)

This is same as:

**unionSets( Find(x), Find(y) )**



## Steps in the Union (x, y)

---

1.  $\text{EqClass}_x = \text{Find}(x)$
2.  $\text{EqClass}_y = \text{Find}(y)$
3.  $\text{EqClass}_{xy} = \text{EqClass}_x \cup \text{EqClass}_y$

union







# A Simple Algorithm to Compute Equivalence Classes

1. Initially, put each element in a set of its own  
i.e.,  $\text{EqClass}_a = \{a\}$ , for every  $a \in S$
  2. FOR EACH element pair  $(a,b)$ :
    1. Check  $[a R b == \text{true}]$
    2. IF  $a R b$  THEN
      1.  $\text{EqClass}_a = \text{Find}(a)$
      2.  $\text{EqClass}_b = \text{Find}(b)$
      3.  $\text{EqClass}_{ab} = \text{EqClass}_a \cup \text{EqClass}_b$
- $\Theta(n^2)$  iterations
- “Union(a,b)”*



# Specification for Union-Find

---

- *Find(x)*
  - Should return the id of the equivalence set that currently contains element x
  
- *Union(a,b)*
  - If a & b are in two different equivalence sets, then Union(a,b) should merge those two sets into one
    - Otherwise, no change



# How to support Union() and Find() operations efficiently?

---

- Approach 1

- *Keep the elements in the form of an array, where:  
 $A[i]$  stores the current set ID for element  $i$*

- Analysis:

- Find() will take  $O(1)$  time
- Union() could take up to  $O(n)$  time
- Therefore a sequence of  $m$  (union and find) operations could take  $O(mn)$  in the worst case
  - This is bad!



# How to support Union() and Find() operations efficiently?

---

- Approach 2
  - *Keep all equivalence sets in separate linked lists:  
1 linked list for every set ID*
  - Analysis:
    - *Union()* now needs only  $O(1)$  time  
(assume doubly linked list)
    - However, *Find()* could take up to  $O(n)$  time
      - Slight improvements are possible (think of Balanced BSTs)
    - A sequence of  $m$  operations takes  $\Omega(m \log n)$
    - Still bad!



# How to support Union() and Find() operations efficiently?

- Approach 3

- *Keep all equivalence sets in separate trees:  
1 tree for every set*
- *Ensure (somehow) that Find() and Union()  
take very little time ( $\ll O(\log n)$ )*

- That is the Union-Find Data Structure!

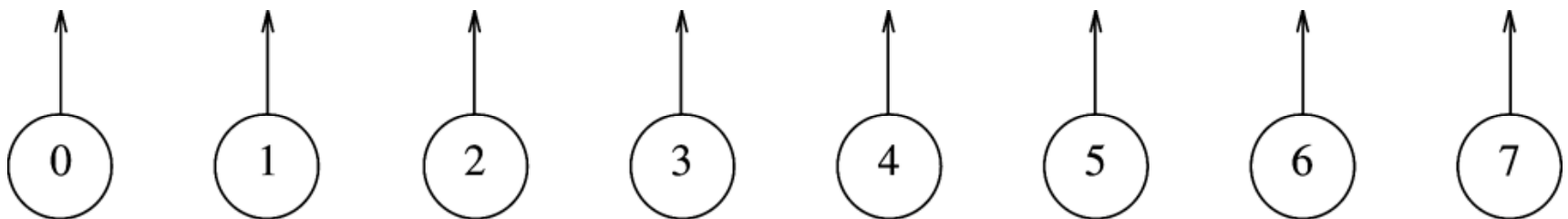
*The Union-Find data structure for  $n$  elements is a forest of  $k$  trees, where  $1 \leq k \leq n$*

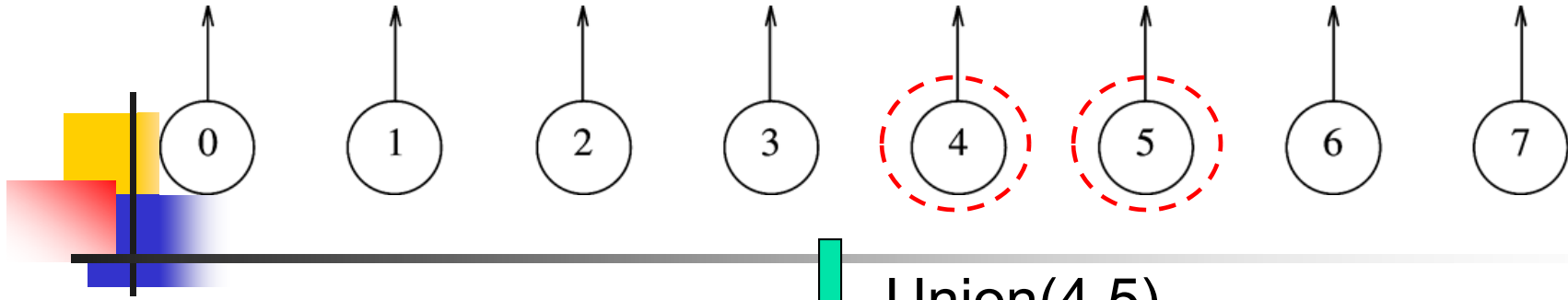


# Initialization

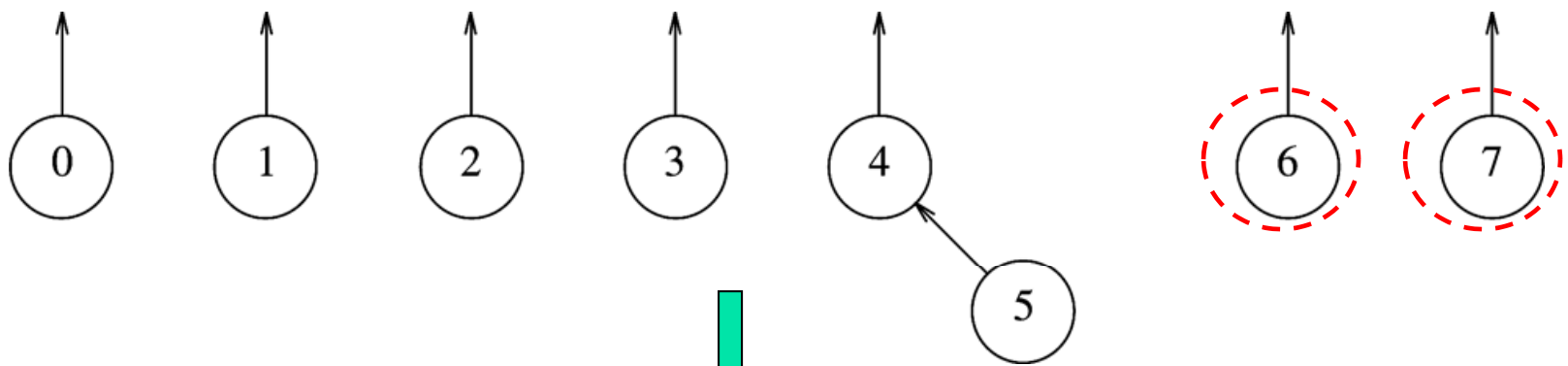
---

- Initially, each element is put in one set of its own
  - Start with  $n$  sets ==  $n$  trees

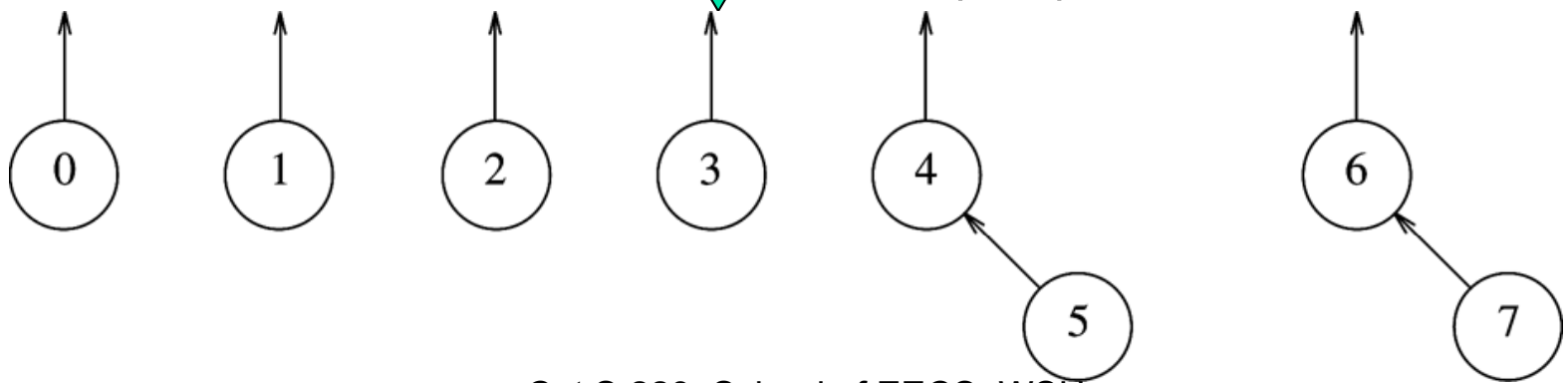


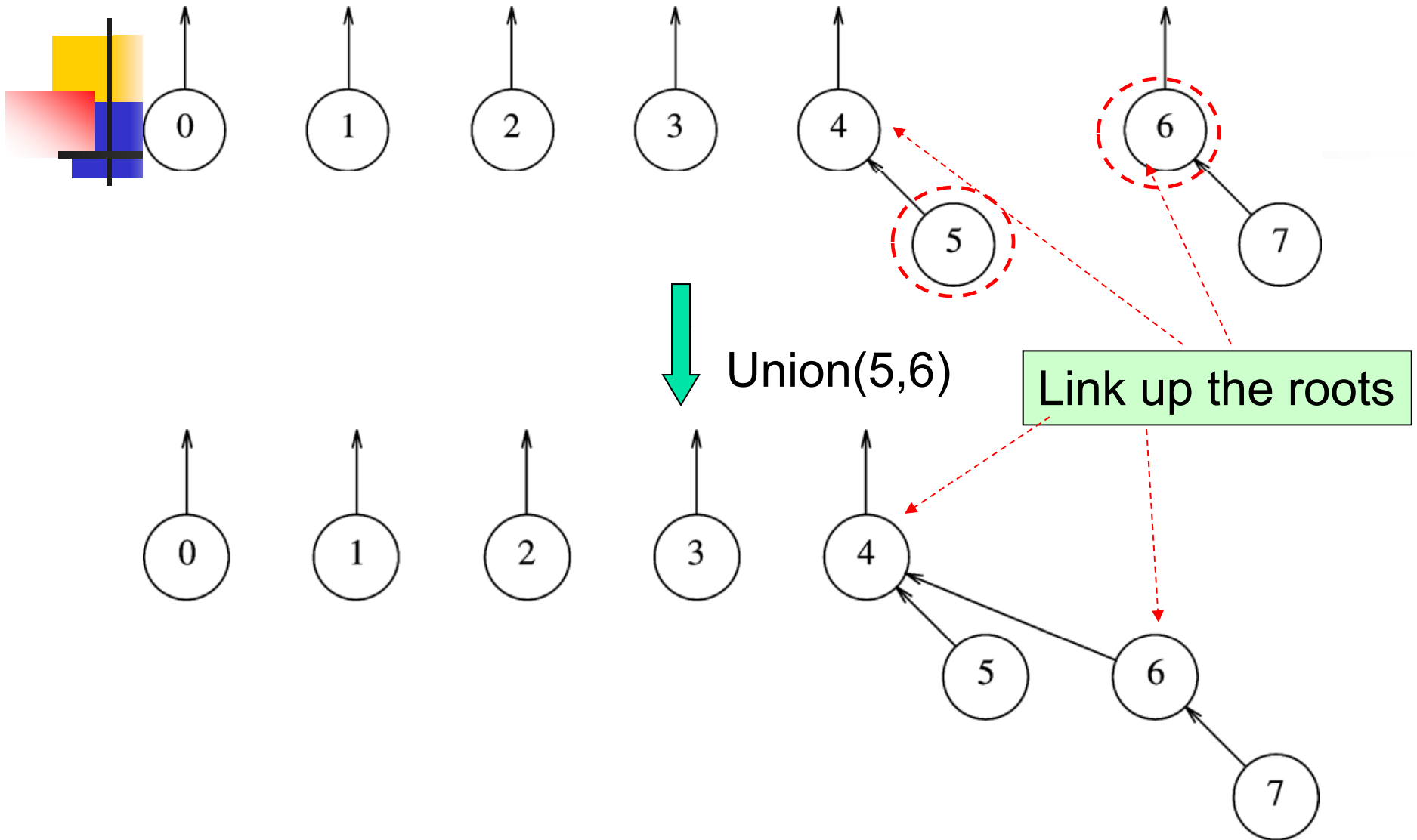


Union(4,5)



Union(6,7)









# The Union-Find Data Structure

---

- **Purpose:** To support two basic operations efficiently
  - *Find* ( $x$ )
  - *Union* ( $x, y$ )
- **Input:** An array of  $n$  elements
- Identify each element by its array index
  - Element label = array index

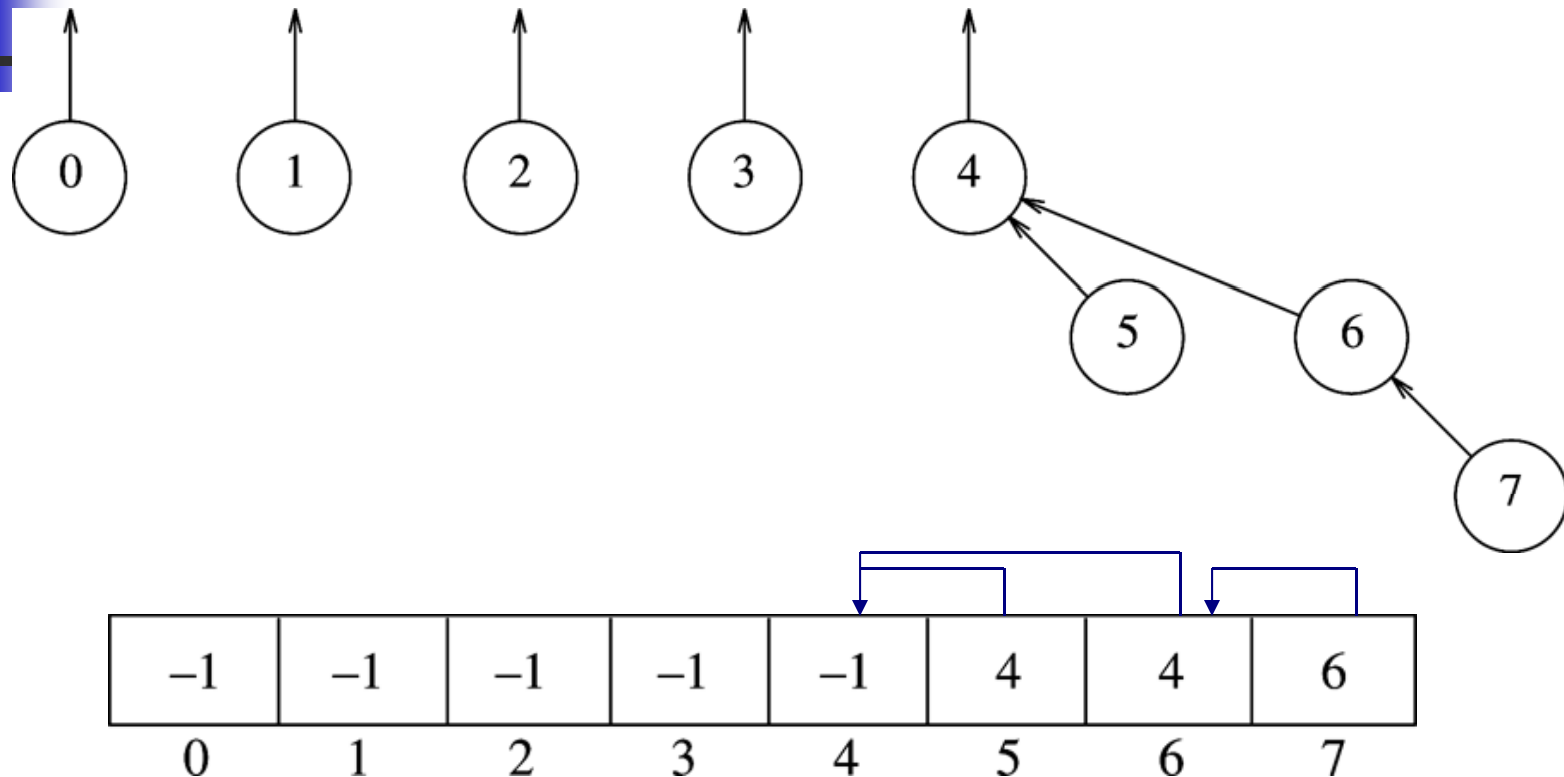


# Union-Find Data Structure

```
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          void union(int x, int y);
10     private:
11         vector<int> s;
12 };
```

Note: This will always be a `vector<int>`, regardless of the data type of your elements. WHY?

# Union-Find D/S: Implementation



- Entry  $s[i]$  points to  $i^{\text{th}}$  parent
- -1 means root

This is WHY  
vector<int>

## **Union-Find: “Simple Version”**

```
1  /**
2   * Construct the disjoint sets object.
3   * numElements is the initial number of disjoint sets.
4   */
5  DisjSets::DisjSets( int numElements ) : s( numElements )
6  {
7      for( int i = 0; i < s.size( ); i++ )
8          s[ i ] = -1;
9  }
```

### **“Simple Find” implementation**

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

### **Union performed arbitrarily**

```
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     s[ root2 ] = root1;
11 }
```

**a & b could be arbitrary  
elements (need not be roots)**

```
void DisjSets::union(int a, int b)
{
    unionSets( find(a), find(b) );
}
```

**This could also be:  
s[root1] = root2  
(both are valid)**



## Analysis of the simple version

---

- Each unionSets() takes only  $O(1)$  in the worst case
- Each Find() could take  $O(n)$  time
  - → Each Union() could also take  $O(n)$  time
- Therefore,  $m$  operations, where  $m \gg n$ , would take  $O(mn)$  in the worst-case

Pretty bad!



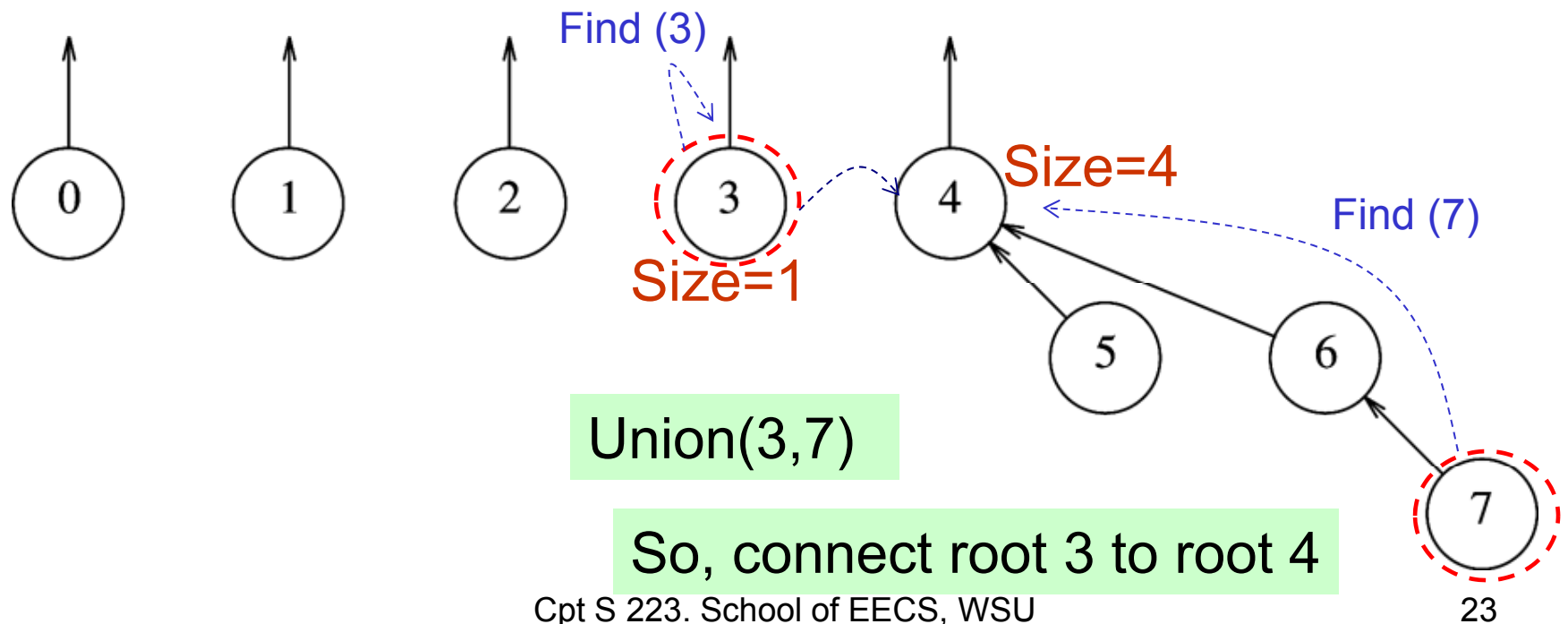
# Smarter Union Algorithms

---

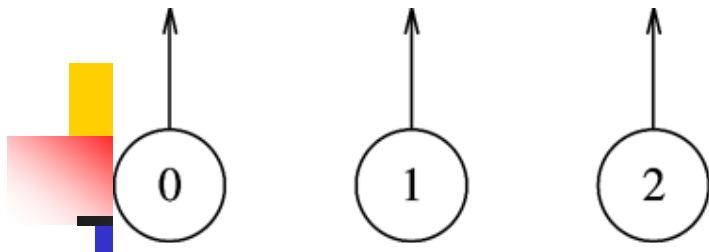
- Problem with the arbitrary root attachment strategy in the simple approach is that:
  - The tree, in the worst-case, could just grow along one long ( $O(n)$ ) path
- Idea: Prevent formation of such long chains
  - => Enforce Union() to happen in a “balanced” way

# Heuristic: *Union-By-Size*

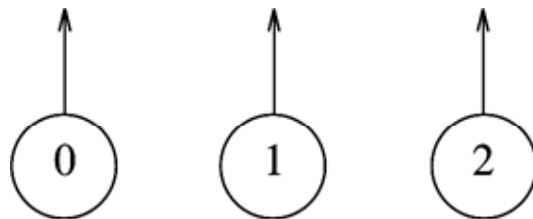
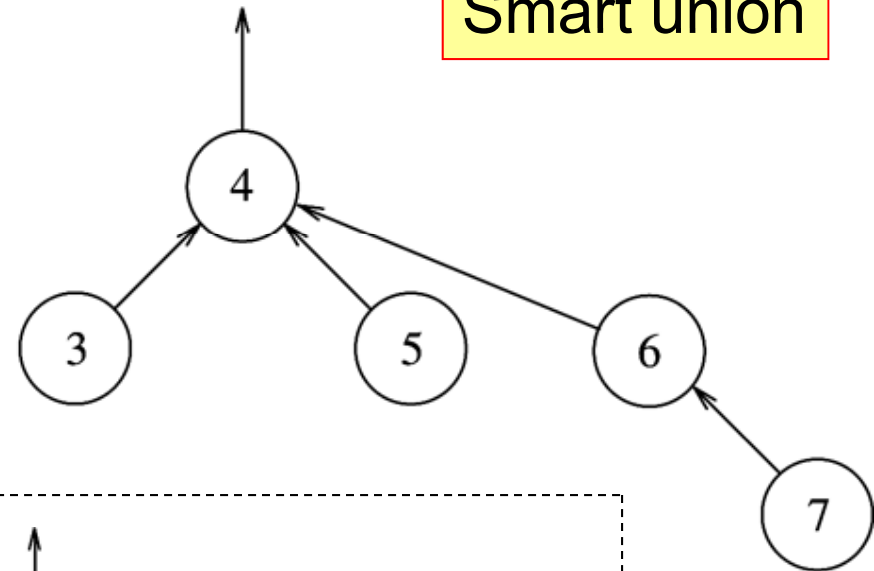
- *Attach the root of the “smaller” tree to the root of the “larger” tree*



## Union-By-Size:

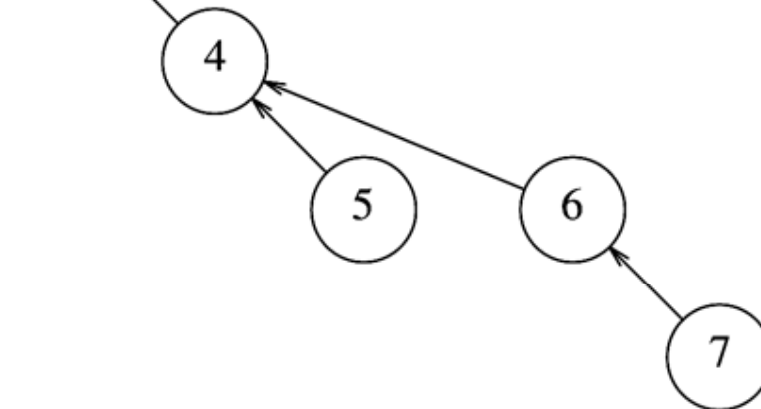


Smart union



Arbitrary Union

An arbitrary union  
could end up  
unbalanced like this:

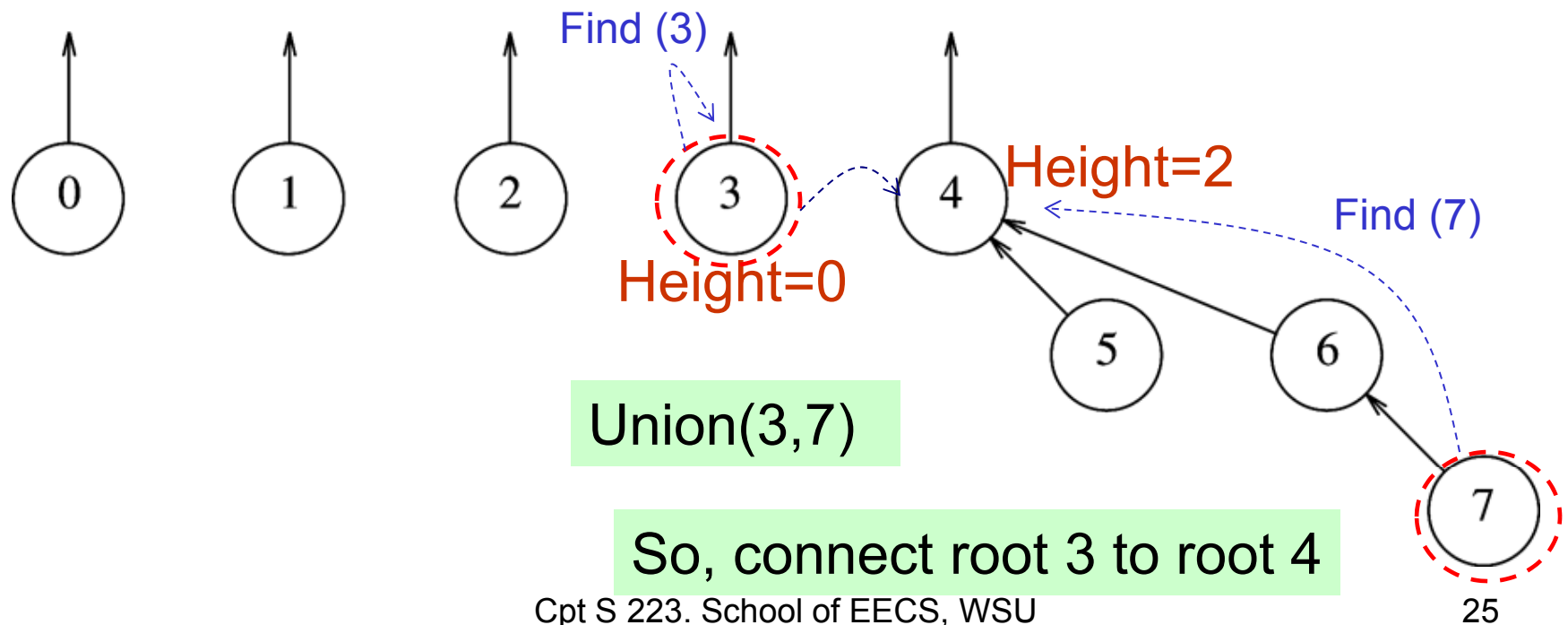




Also known as  
“*Union-By-Rank*”

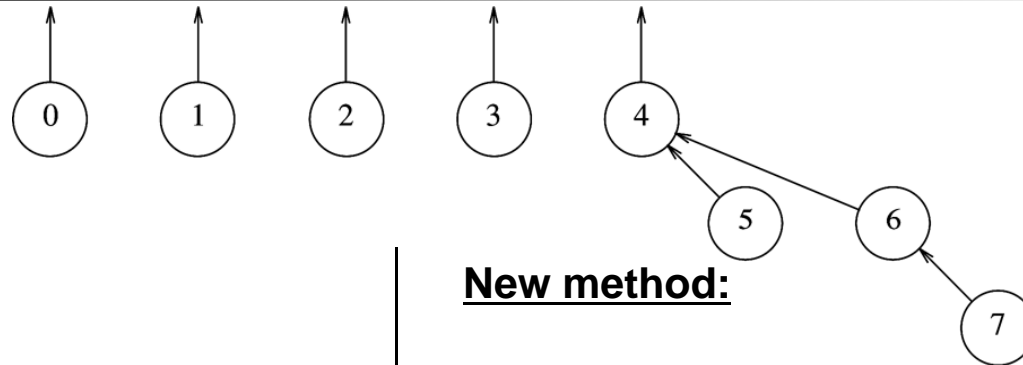
## Another Heuristic: *Union-By-Height*

- *Attach the root of the “shallower” tree to the root of the “deeper” tree*



Let us assume union-by-rank first

# How to implement smart union?



Old method:

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

But where will you keep track of the heights?

- $s[i]$  = parent of  $i$
- $S[i] = -1$ , means root

New method:

What is the problem if you store the height value directly?

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

- instead of roots storing -1, let them store a value that is equal to:  $-1 - (\text{tree height})$



# New code for union by rank?

---

```
void DisjSets::unionSets(int root1,int root2) {  
  
    // first compare heights  
  
    // link up shorter tree as child of taller tree  
    // if equal height, make arbitrary choice  
  
    // then increment height of new merged tree if height  
    // has changed – will happen if merging two equal  
    // height trees  
  
}
```



# New code for union by rank?

---

```
void DisjSets::unionSets(int root1,int root2) {  
    assert(s[root1]<0);  
    assert(s[root2]<0);  
    if(s[root1]<s[root2])        s[root2]=root1;  
    if(s[root2]<s[root1])        s[root1]=root2;  
    if(s[root1]==s[root2]) {  
        s[root1]=root2;  
        s[root2]--;  
    }  
}
```

Note: All nodes, except root, store parent id.  
The root stores a value = negative(height) -1

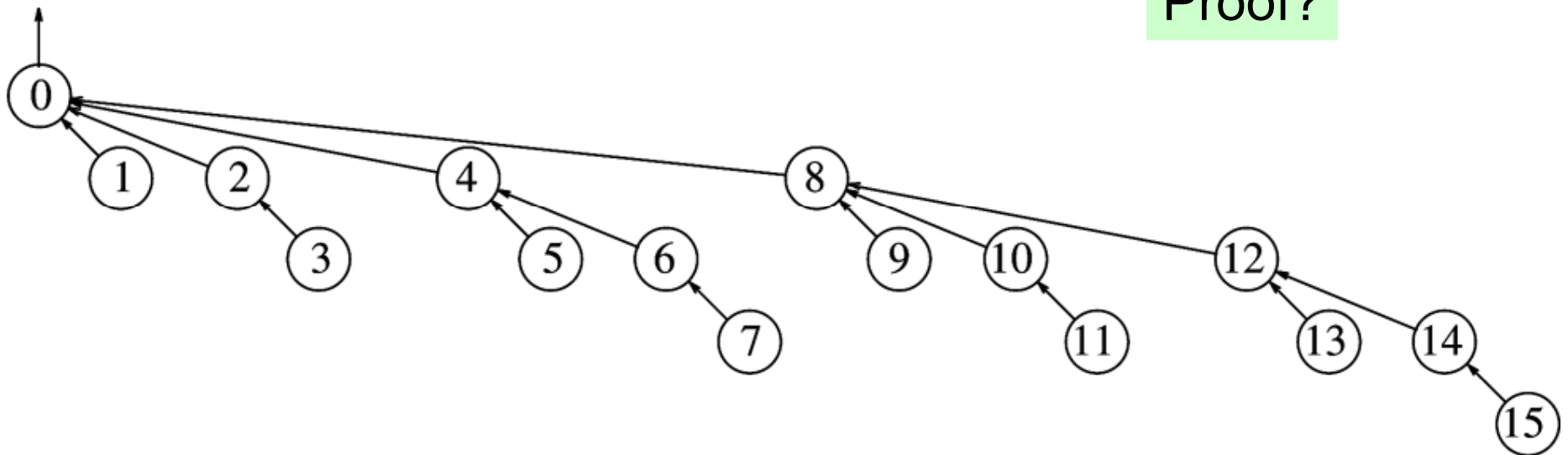
## Code for Union-By-Rank

```
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 }
```

Similar code for  
union-by-size

# How Good Are These Two Smart Union Heuristics?

## ■ Worst-case tree



Maximum depth restricted to  $O(\log n)$



# Analysis: Smart Union Heuristics

---

- For smart union (by rank or by size):
  - Find() takes  $O(\log n)$ ;
    - $\implies$  union() takes  $O(\log n)$ ;
  - unionSets() takes  $O(1)$  time
- For  $m$  operations:  $O(m \log n)$  run-time
- Can it be better?
  - What is still causing the  $(\log n)$  factor is the distance of the root from the nodes
  - Idea: Get the nodes as close as possible to the root
    - Path Compression!**

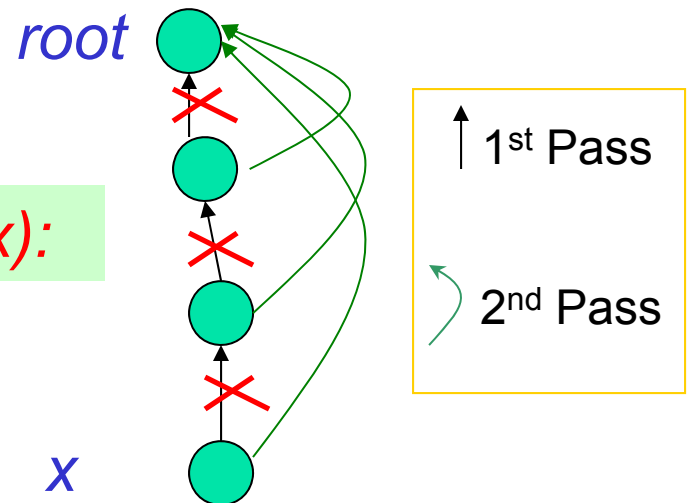
# Path Compression Heuristic

- During find(x) operation:
  - Update all the nodes along the path from x to the root point directly to the root
  - A two-pass algorithm

How will this help?

Any future calls to find on x or its ancestors will return in constant time!

*find(x):*







## New code for find() using path compression?

---

```
void DisjSets::find(int x) {
```

```
?
```

```
}
```



# New code for find() using path compression?

---

```
int DisjSets::find(int x) {  
  
    // if x is root, then just return x  
    if(s[x]<0) return x;  
  
    // otherwise simply call find recursively, but..  
    //     make sure you store the return value (root index)  
    //           to update s[x], for path compression  
  
    return s[x]=find(s[x]);  
}
```



# Path Compression: Code

```
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 }
```

It can be proven that path compression alone ensures that  $\text{find}(x)$  can be achieved in  $O(\log n)$

Spot the difference from old  $\text{find}()$  code!

## Union-by-Rank & Path-Compression: Code

*Init()*

```
1  /**
2   * Construct the disjoint sets object.
3   * numElements is the initial number of disjoint sets.
4   */
5  DisjSets::DisjSets( int numElements ) : s( numElements )
6  {
7      for( int i = 0; i < s.size( ); i++ )
8          s[ i ] = -1;
9  }
```

### *Union()*

```
void DisjSets::union(int a, int b)
{
    unionSets( find(a), find(b) );
}
```

*unionSets()*

**Smart union**

```
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 }
```

*Find()*

**Smart find**

```
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 }
```


Amortized complexity for m operations:  
 $O(m \text{ Inv. Ackerman}(m,n)) = O(m \log^* n)$



# Heuristics & their Gains

	Worst-case run-time for m operations
Arbitrary Union, Simple Find	$O(m n)$
Union-by-size, Simple Find	$O(m \log n)$
Union-by-rank, Simple Find	$O(m \log n)$
Arbitrary Union, Path compression Find	$O(m \log n)$
<b>Union-by-rank, Path compression Find</b>	$O(m \text{ Inv.Ackermann}(m,n))$ $= O(m \log^* n)$

Extremely slow  
Growing function





# What is Inverse Ackermann Function?

---

- $A(1, j) = 2^j$  for  $j \geq 1$
- $A(i, 1) = A(i-1, 2)$  for  $i \geq 2$
- $A(i, j) = A(i-1, A(i, j-1))$  for  $i, j \geq 2$
- $\text{InvAck}(m, n) = \min\{i \mid A(i, \text{floor}(m/n)) > \log N\}$
- $\text{InvAck}(m, n) = O(\log^* n)$  (pronounced “log star n”)
  - Even Slower!
  - A very slow function



# How Slow is Inverse Ackermann Function?

---

- What is  $\log^*n$ ?

← A very slow function

- $\log^*n = \log \log \log \log \dots n$ 
  - How many times we have to repeatedly take log on  $n$  to make the value to 1?
- $\log^*65536=4$ , but  $\log^*2^{65536}=5$



# Some Applications

---



# A Naïve Algorithm for Equivalence Class Computation

1. Initially, put each element in a set of its own

i.e.,  $\text{EqClass}_a = \{a\}$ , for every  $a \in S$

$\Theta(n^2)$   
iterations

2. FOR EACH element pair (a,b):

1. Check  $[a R b == \text{true}]$

2. IF  $a R b$  THEN

1.  $\text{EqClass}_a = \text{Find}(a)$

2.  $\text{EqClass}_b = \text{Find}(b)$

3.  $\text{EqClass}_{ab} = \text{EqClass}_a \cup \text{EqClass}_b$

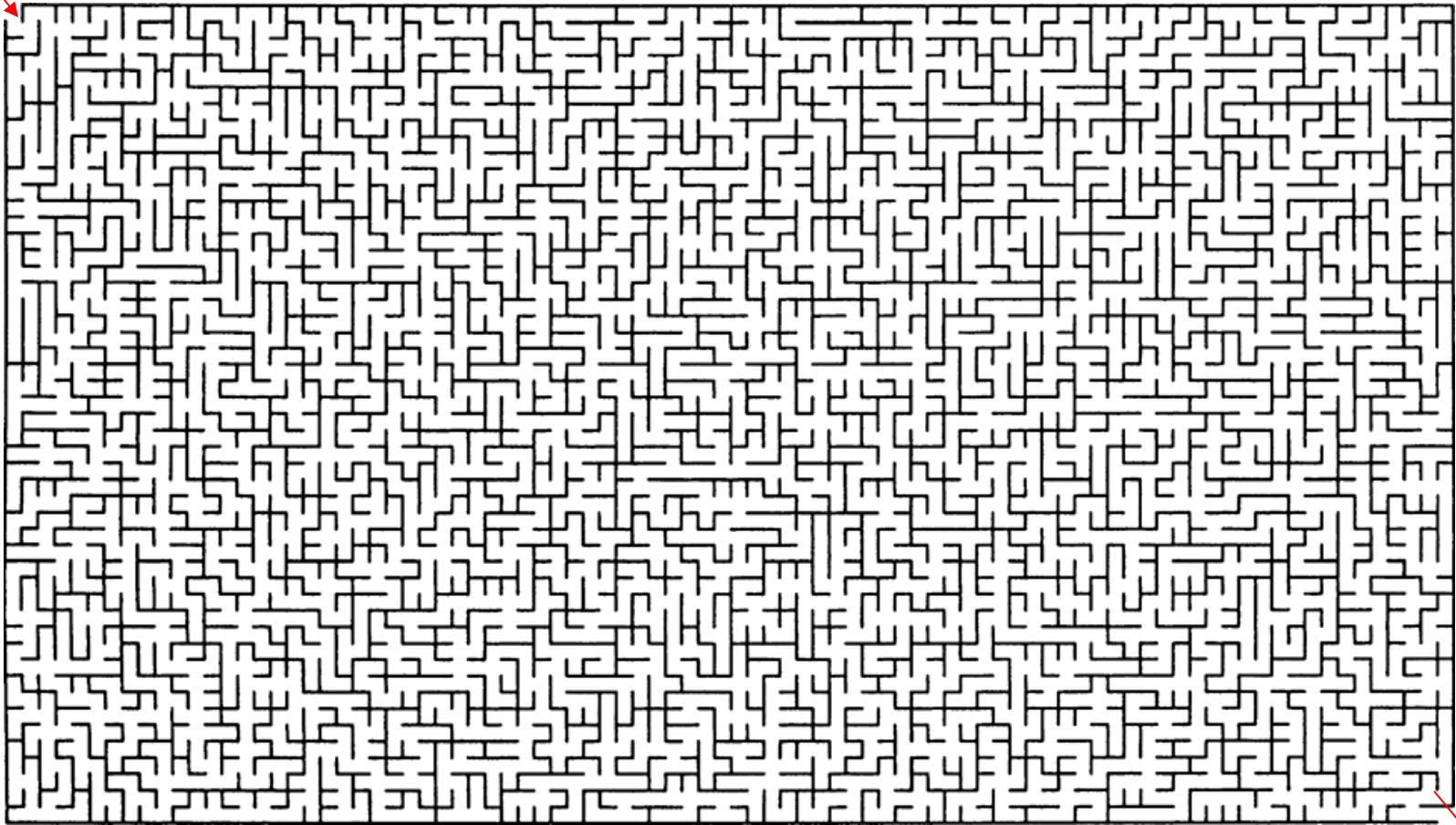
$O(\log^* n)$   
amortized

Run-time using union-find:  $O(n^2 \log^* n)$

Better solutions using other data structures/techniques could exist depending on the application



# An Application: Maze



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

### Strategy:

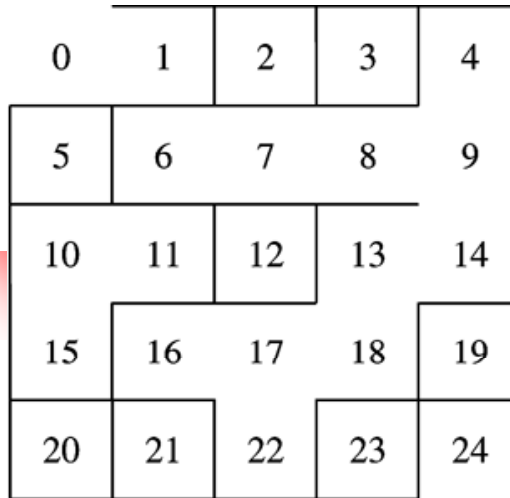
- As you find cells that are connected, collapse them into equivalent set
- If no more collapses are possible, examine if the Entrance cell and the Exit cell are in the same set
  - If so => we have a solution
  - O/w => no solutions exists

{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

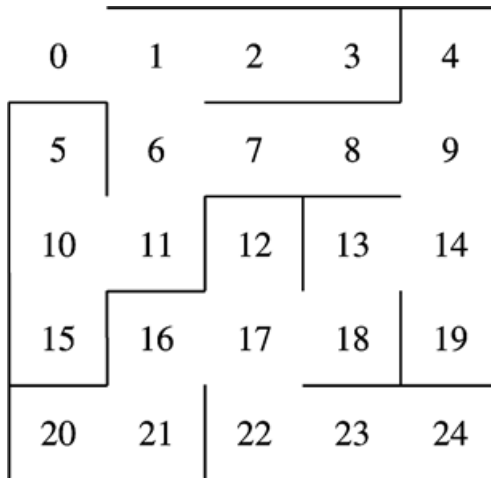
{0, 1} {2} {3} {4, 6, 7, 8, 9, 13, 14} {5} {10, 11, 15} {12} {16, 17, 18, 22} {19} {20} {21} {23} {24}



### Strategy:

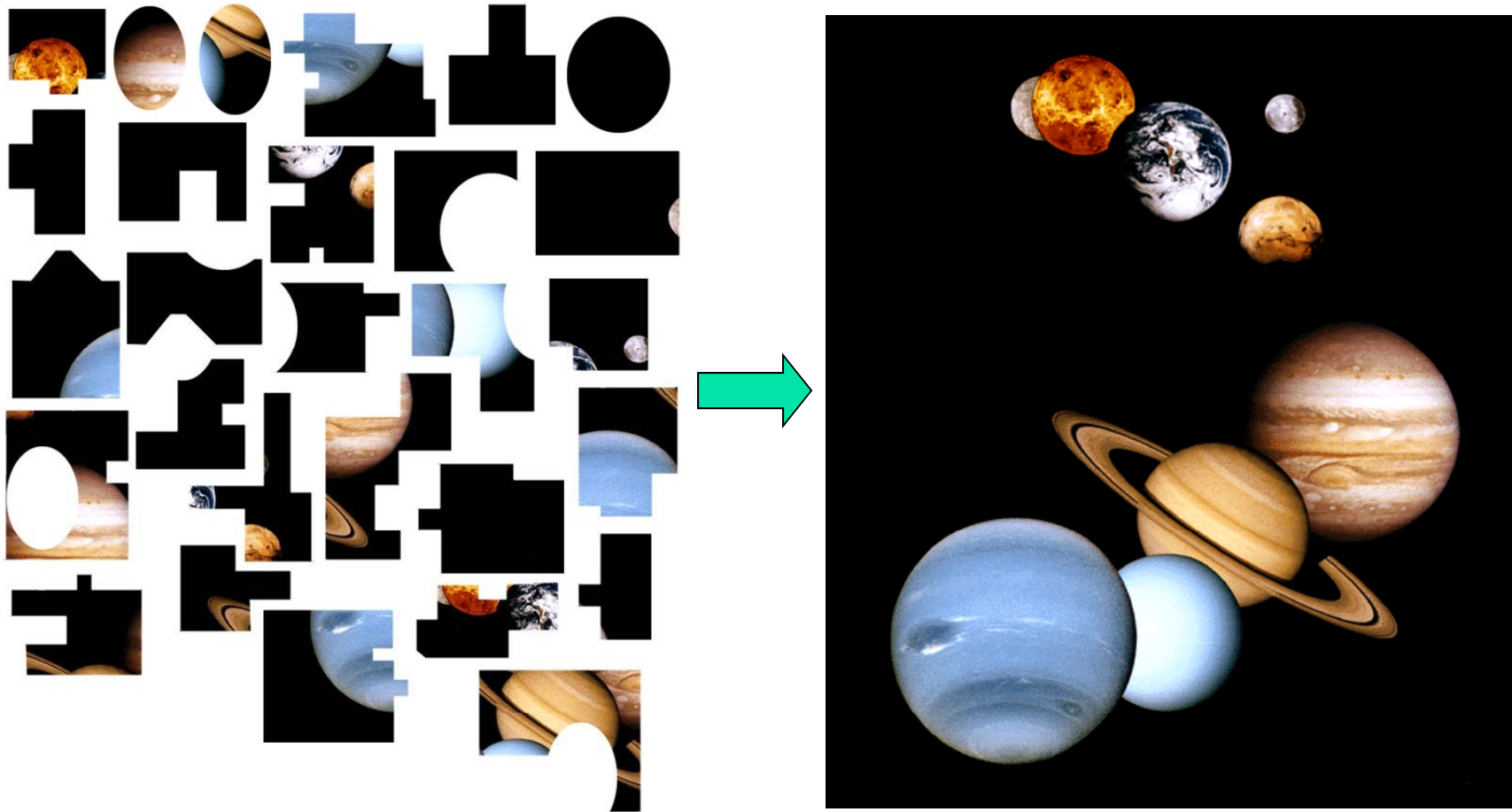
- As you find cells that are connected, collapse them into equivalent set
- If no more collapses are possible, examine if the Entrance cell and the Exit cell are in the same set
  - If so => we have a solution
  - O/w => no solutions exists

{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}



{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}

# Another Application: Assembling Multiple Jigsaw Puzzles at once



Merging Criterion: Visual & Geometric Alignment

Picture Source: <http://ssed.gsfc.nasa.gov/lepedu/jigsaw.html>



# Summary

---

- Union Find data structure
  - Simple & elegant
  - Complicated analysis
- Great for disjoint set operations
  - Union & Find
  - In general, great for applications with a need for “clustering”