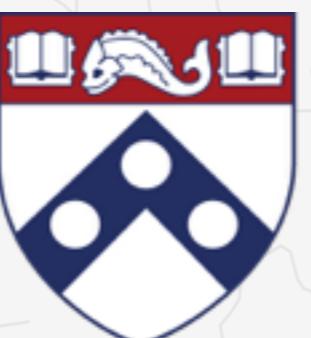


CIS 121, FALL 2015

DATA STRUCTURES AND ALGORITHMS

AKA PROGRAMMING LANGUAGES & TECHNIQUES II

CHRIS CALLISON-BURCH



Penn
UNIVERSITY of PENNSYLVANIA

CIS 121 course overview

What is CIS 121?

- Third course in the intro sequence CIS 120, 160, 121
- Programming and problem solving, with applications.
- **Algorithm:** method for solving a problem.
- **Data structure:** method to store information.

topic	data structures and algorithms
data types	stack, queue, bag, union-find, priority queue
sorting	quicksort, mergesort, heapsort, radix sorts
searching	BST, red-black BST, hash table
graphs	BFS, DFS, Prim, Kruskal, Dijkstra
strings	KMP, regular expressions, tries, data compression
advanced	B-tree, k-d tree, suffix array, maxflow

Why study algorithms?

Their impact is broad and far-reaching.

Internet. Web search, packet routing, distributed file sharing, ...

Biology. Human genome project, protein folding, ...

Computers. Circuit layout, file system, compilers, ...

Computer graphics. Movies, video games, virtual reality, ...

Security. Cell phones, e-commerce, voting machines, ...

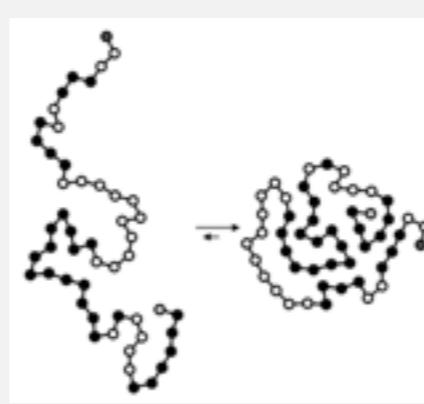
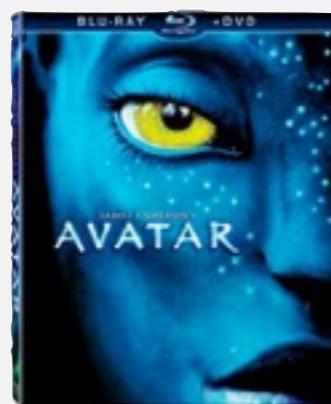
Multimedia. MP3, JPG, DivX, HDTV, face recognition, ...

Social networks. Recommendations, news feeds, advertisements, ...

Physics. N-body simulation, particle collision simulation, ...

:

Google
YAHOO!
bing



Why study algorithms?

Their impact is broad and far-reaching.

Mysterious algorithm was 4% of trading activity last week

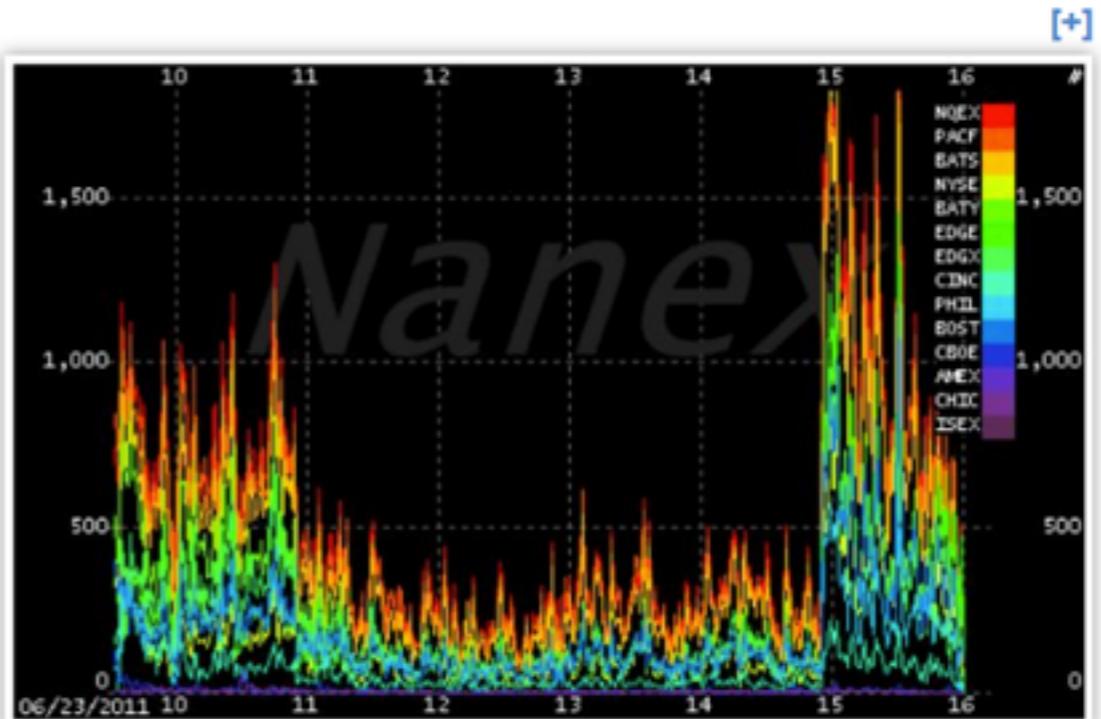
October 11, 2012

A single mysterious computer program that placed orders — and then subsequently canceled them — made up 4 percent of all quote traffic in the U.S. stock market last week, according to the top tracker of [high-frequency trading](#) activity.

The motive of the algorithm is still unclear, [CNBC](#) reports.

The program placed orders in 25-millisecond bursts involving about 500 stocks, according to Nanex, a market data firm. The algorithm never executed a single trade, and it abruptly ended at about 10:30 a.m. ET Friday.

"My guess is that the algo was testing the market, as high-frequency frequently does," says Jon Najarian, co-founder of TradeMonster.com. "As soon as they add bandwidth, the HFT crowd sees how quickly they can top out to create latency." ([Read More: Unclear What Caused Kraft Spike: Nanex Founder.](#))



Generic high frequency trading chart (credit: Nanex)

Why study algorithms?

For intellectual stimulation.

“For me, great algorithms are the poetry of computation. Just like verse, they can be terse, allusive, dense, and even mysterious.

But once unlocked, they cast a brilliant new light on some aspect of computing. ” — Francis Sullivan

FROM THE EDITORS

THE JOY OF ALGORITHMS

Francis Sullivan, Associate Editor-in-Chief

THE THEME OF THIS FIRST-OF-THE-CENTURY ISSUE OF COMPUTING IN SCIENCE & ENGINEERING IS ALGORITHMS. IN FACT, WE WERE BOLD ENOUGH—AND PERHAPS FOOLISH ENOUGH—to call the 10 examples we've selected “THE TOP 10 ALGORITHMS OF THE CENTURY.”

Computational algorithms are probably as old as civilization. Sumerian cuneiform, one of the most ancient written records, contains what may be the first algorithm, for calculating the area of a trapezoid. The algorithm was used for 4,000 years in his whole life. He wasn't joking, because he was referring to the 15 minutes during which he'd sketched out a fundamental operation of geometry. The algorithm was a simple one of thought and investigation in a sink cost that night or might not have paid off.

The last century has cracked many hard problems since 1 January 1900, but we are posing some even harder ones on to the next century. In spite of a lot of good work, a lot of useful predictions, and a lot of progress, there is still a very large mass of data that is still almost untouched. There are still very big challenges coming from more “traditional” information sciences. For example, recall one late-night session on the Maryland shore where the question was “What's the area of a crab? After all, they don't look very appealing.” After the usual quibbling about what “area” means, the question was asked: “What must be the right answer—namely, ‘A very hungry person first ate a crab.’”

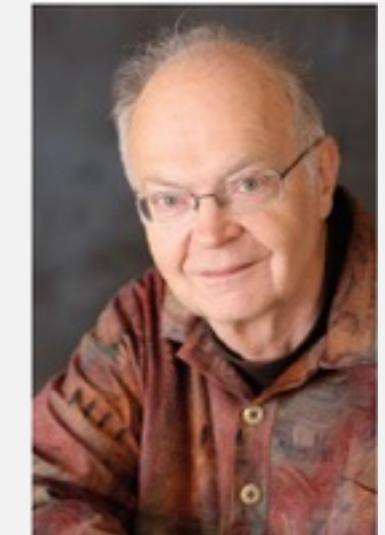
The problem here is that the number of “computations” to “compute” creates its own necessity. Our need for powerful methods always exceeds their availability. Each significant computation is a challenge to the computer designer. As the problem gets larger, computation to be done. New algorithms are an attempt to bridge the gap between the demand for cycles and the available time. The problem is that the time constant is increasing at Moore's Law factor of two every 18 months. In effect, Moore's Law changes the constant in front of the estimate of running time. The result is that the time constant is decreasing. The cost does not come along every 1.5 years, but when they do, they can change the expansion of the computation!

For the last century, the mystery of computation is over. It is there an analog for things such as huge, multidisciplinary projects? I suspect that the answer is no. There are no comparable methods for solving specific cases of “impossible” problems. Instances of NP-complete problems crop up in almost every field of science and engineering. Are there efficient ways to attack them?

I suspect that in the 21st century, there will be a race for an understanding of the basic nature of computation, of computational theory. Questions already arising from quantum computing and problems associated with the generation and storage of data, the nature of computation, and the interplay together theories of computing, logic, and the nature of the physical world.

The next century is not going to be very fruitful for us, but it is not going to be dull either. ■

“An algorithm must be seen to be believed. ” — Donald Knuth

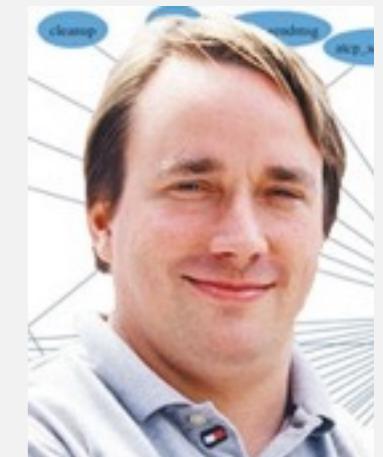


Why study algorithms?

To become a proficient programmer.

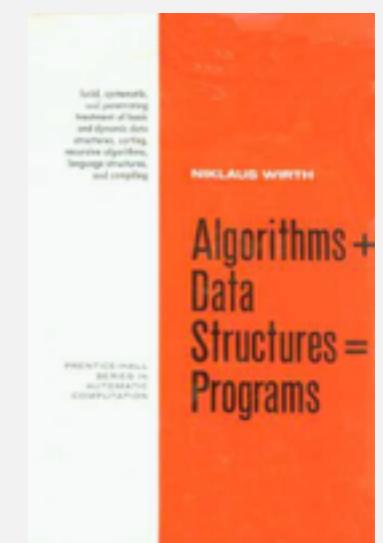
“I will, in fact, claim that the difference between a bad programmer and a good one is whether whether the programmer considers code or data structures more important. Bad programmers worry about the code. Good programmers worry about data structures and their relationships.”

— Linus Torvalds (creator of Linux)



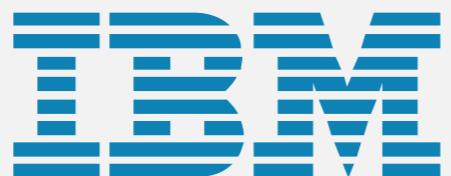
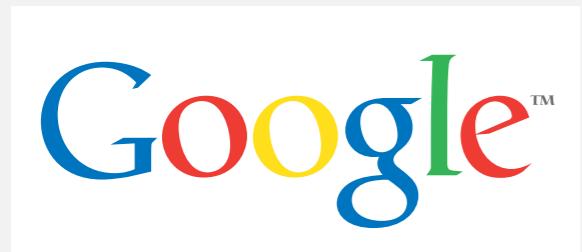
“Algorithms + Data Structures = Programs.”

— Niklaus Wirth



Why study algorithms?

For the interview.



Why study algorithms?

They may unlock the secrets of life and of the universe.

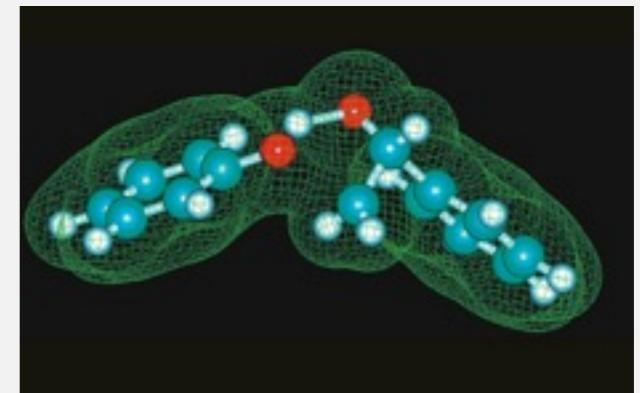
“ Computer models mirroring real life have become crucial for most advances made in chemistry today.... Today the computer is just as important a tool for chemists as the test tube. ”

— Royal Swedish Academy of Sciences

(Nobel Prize in Chemistry 2013)



Martin Karplus, Michael Levitt, and Arieh Warshel



Why study algorithms?

- Their impact is broad and far-reaching.
- For intellectual stimulation.
- To become a proficient programmer.
- To pass your job interviews.
- They may unlock the secrets of life and of the universe.

Why study anything else?



Coursework and grading

Weekly assignments and term project. 55%

- Mix of programming and written assignments
- Collaboration/lateness policies: see web.

Exams. 10% + 10% + 20%

- Midterm 1 (in class on Thursday, October 1).
- Midterm 2 (in class on Tuesday, November 3)
- Final (9am-11am on Tuesday, December 15).

Recitations. 5%

- Attendance is mandatory

Resources (textbook)

Required reading. Algorithms 4th edition by R. Sedgewick and K. Wayne.
This course closely follows the textbook and uses their lectures.



Available in hardcover and Kindle.

- Online: \$60/\$55 to buy, \$18 to rent

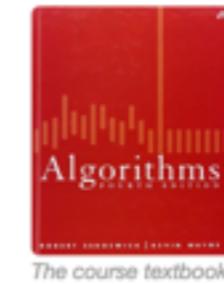
Resources (web)

Course content

- Course info
- Lecture schedule
- Homework assignments
- Lab handouts
- Contact info and office hours
- Pointers to other information

CIS 121 - Data Structures and Algorithms

Lectures Labs Homework Final Project Staff Office Hours Tutoring Resources



Course number:
CIS 121
Instructor:
[Chris Callison-Burch](#)
Teaching Assistants:
[Course Staff](#)
Discussion Forum:
[Piazza](#)
Time and place:
Lectures are Tuesdays and Thursdays at 10:30 in Towne 100
Prerequisites:
[CIS 120 and CIS 160](#)
Textbook:
The lectures and readings will closely follow [Algorithms by Sedgewick and Wayne, 4th edition](#). The textbook has its own web site with some free materials, but you should also buy the full textbook.

<http://www.seas.upenn.edu/~cis121/>

Booksite

- Brief summary of content
- Download code from book

The image shows the front cover of the textbook 'Algorithms, 4th Edition'. The cover is red with the title 'ALGORITHMS, 4TH EDITION' in white. Below the title, it says 'essential information that every serious programmer needs to know about algorithms and data structures'. At the bottom, it says 'ROBERT SEDGEWICK KEVIN WAYNE'.

Textbook. The textbook *Algorithms, 4th Edition* by Robert Sedgewick and Kevin Wayne [[Amazon](#) · [Addison-Wesley](#)] surveys the most important algorithms and data structures in use today. The textbook is organized into six chapters:

- *Chapter 1: Fundamentals* introduces a scientific and engineering basis for comparing algorithms and making predictions. It also includes our programming model.
- *Chapter 2: Sorting* considers several classic sorting algorithms, including insertion sort, mergesort, and quicksort. It also includes a binary heap implementation of a priority queue.
- *Chapter 3: Searching* describes several classic symbol table implementations, including binary search trees, red-black trees, and hash tables.

<http://algs4.cs.princeton.edu>

Where to get help?

Piazza discussion forum.

- Low latency, low bandwidth.
- Mark solution-revealing questions as private.

piazza

<https://piazza.com/upenn/fall2015/cis121/home>

Office hours.

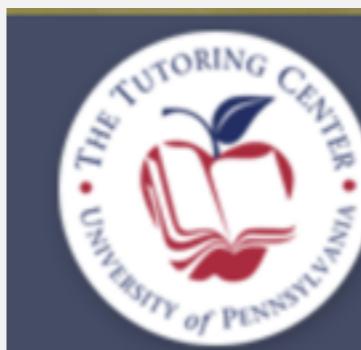
- High bandwidth, high latency.
- See web for schedule.



<http://www.seas.upenn.edu/~cis121/>

Tutoring.

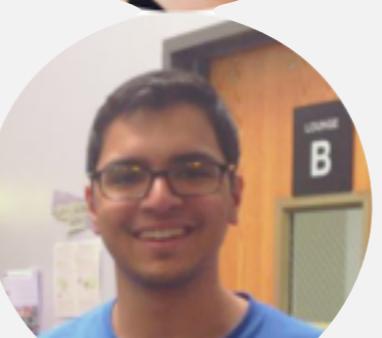
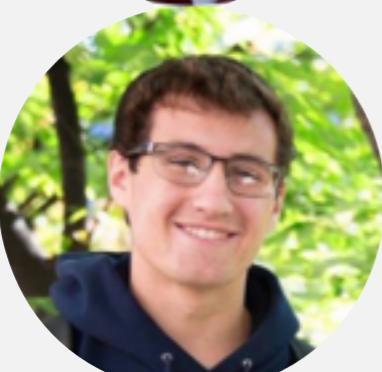
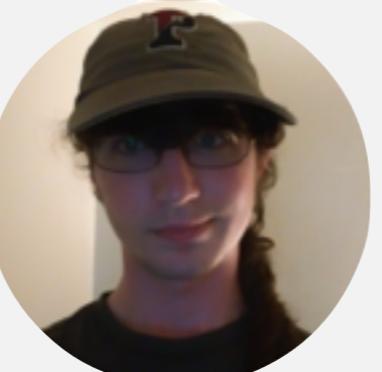
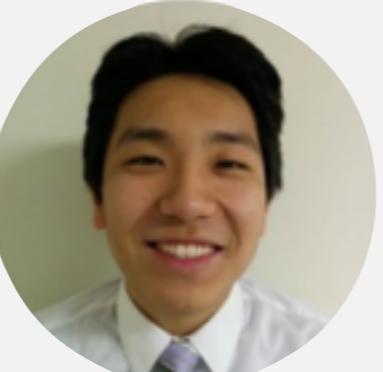
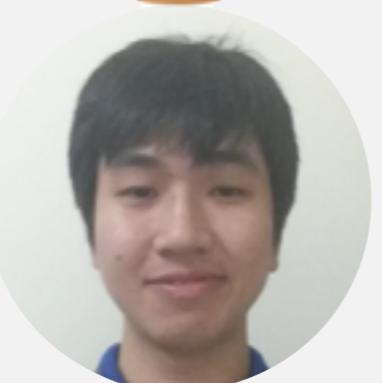
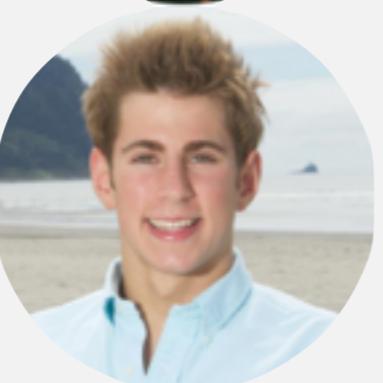
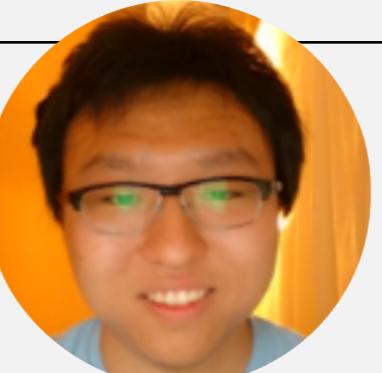
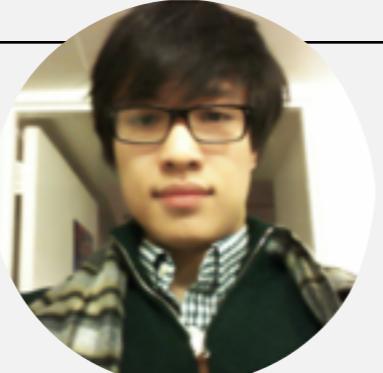
- The Tutoring Center



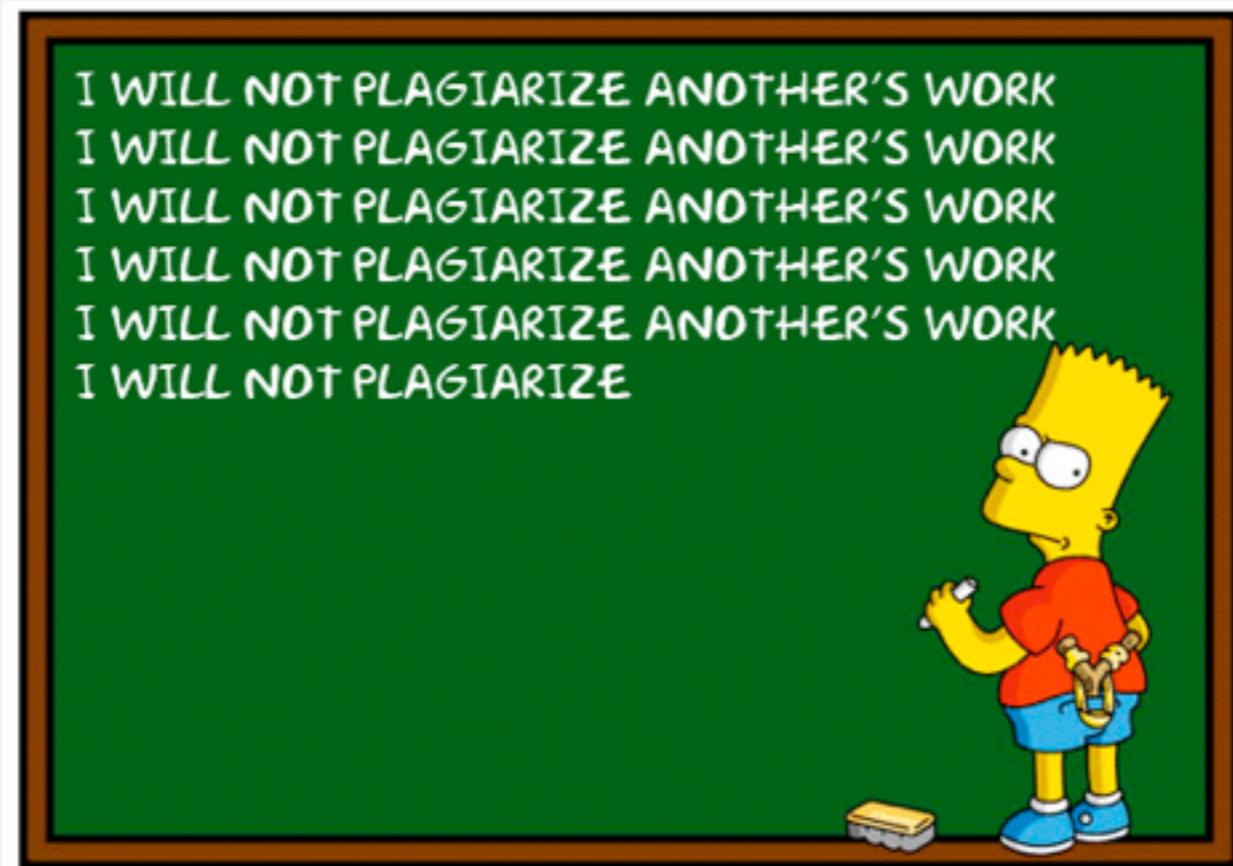
The Tutoring Center

<http://www.vpul.upenn.edu/tutoring/>

Course Staff



Where not to get help?



<http://world.edu/academic-plagiarism>

http://www.upenn.edu/academicintegrity/ai_codeofacademicintegrity.html

How to succeed in CIS 121

Attend lectures and read the textbook

Start homework assignments early

- Do not wait until the last day

Make use of your TAs and instructors

- Attend recitations
- Go to office hours (fine to come even if you don't have a specific question)

Do not cheat

What's ahead?

Lecture 1. [today] Union find.

Next Week. [Tuesday and Thursday] CCB will be traveling for DARPA meeting. Lectures by Rajiv Gandhi. Review of CIS 160.



Recitations start on Monday

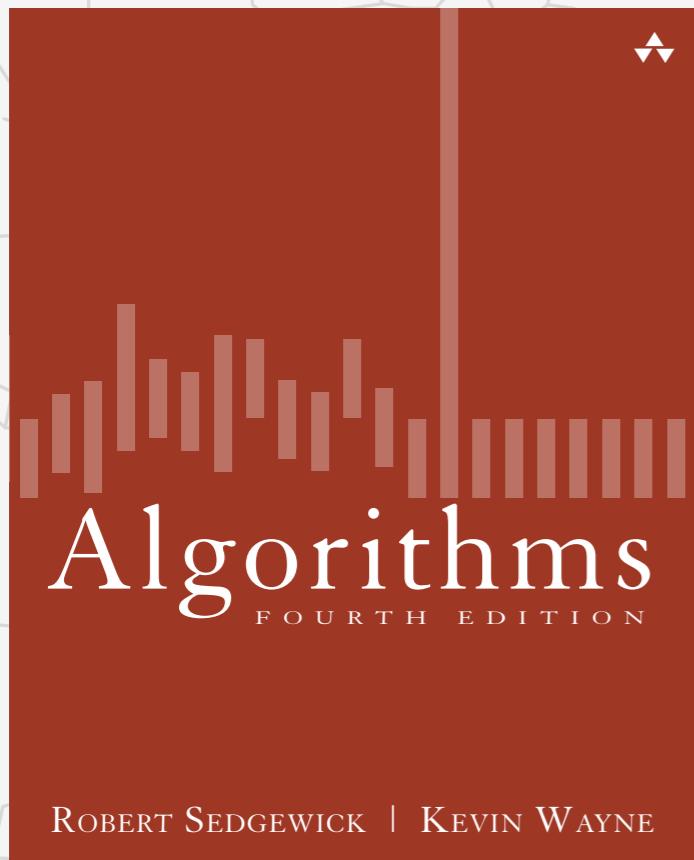
Homework 1. Released on Tuesday. Due via electronic submission before class on September 8, a week from Tuesday.

protip: start early

Not registered? Sign up for the waitlist on the web site, go to any recitation next week.

Algorithms

ROBERT SEDGEWICK | KEVIN WAYNE



ROBERT SEDGEWICK | KEVIN WAYNE

<http://algs4.cs.princeton.edu>

1.5 UNION-FIND

- ▶ *dynamic connectivity*
- ▶ *quick find*
- ▶ *quick union*
- ▶ *improvements*
- ▶ *applications*

Theme of today's lecture (and this course)

Steps to developing a usable algorithm.

- Model the problem.
- Find an algorithm to solve it.
- Fast enough? Fits in memory?
- If not, figure out why not.
- Find a way to address the problem.
- Iterate until satisfied.

The scientific method.

Mathematical analysis.

Algorithms

ROBERT SEDGEWICK | KEVIN WAYNE

<http://algs4.cs.princeton.edu>

1.5 UNION-FIND

- ▶ *dynamic connectivity*
- ▶ *quick find*
- ▶ *quick union*
- ▶ *improvements*
- ▶ *applications*

Dynamic connectivity problem

Given a set of N objects, support two operation:

- Connect two objects.
- Is there a path connecting the two objects?

connect 4 and 3

connect 3 and 8

connect 6 and 5

connect 9 and 4

connect 2 and 1

are 0 and 7 connected? ✗

are 8 and 9 connected? ✓

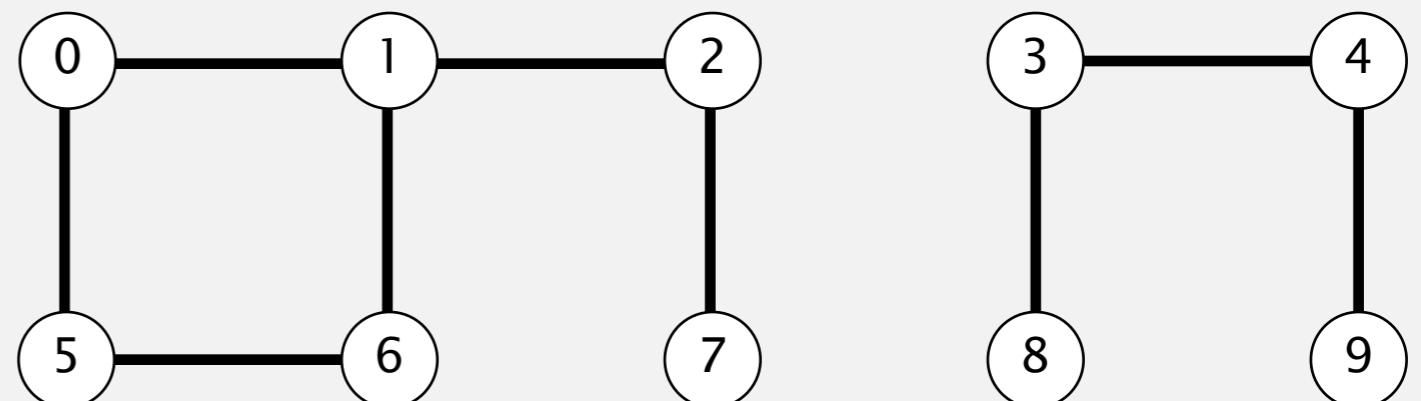
connect 5 and 0

connect 7 and 2

connect 6 and 1

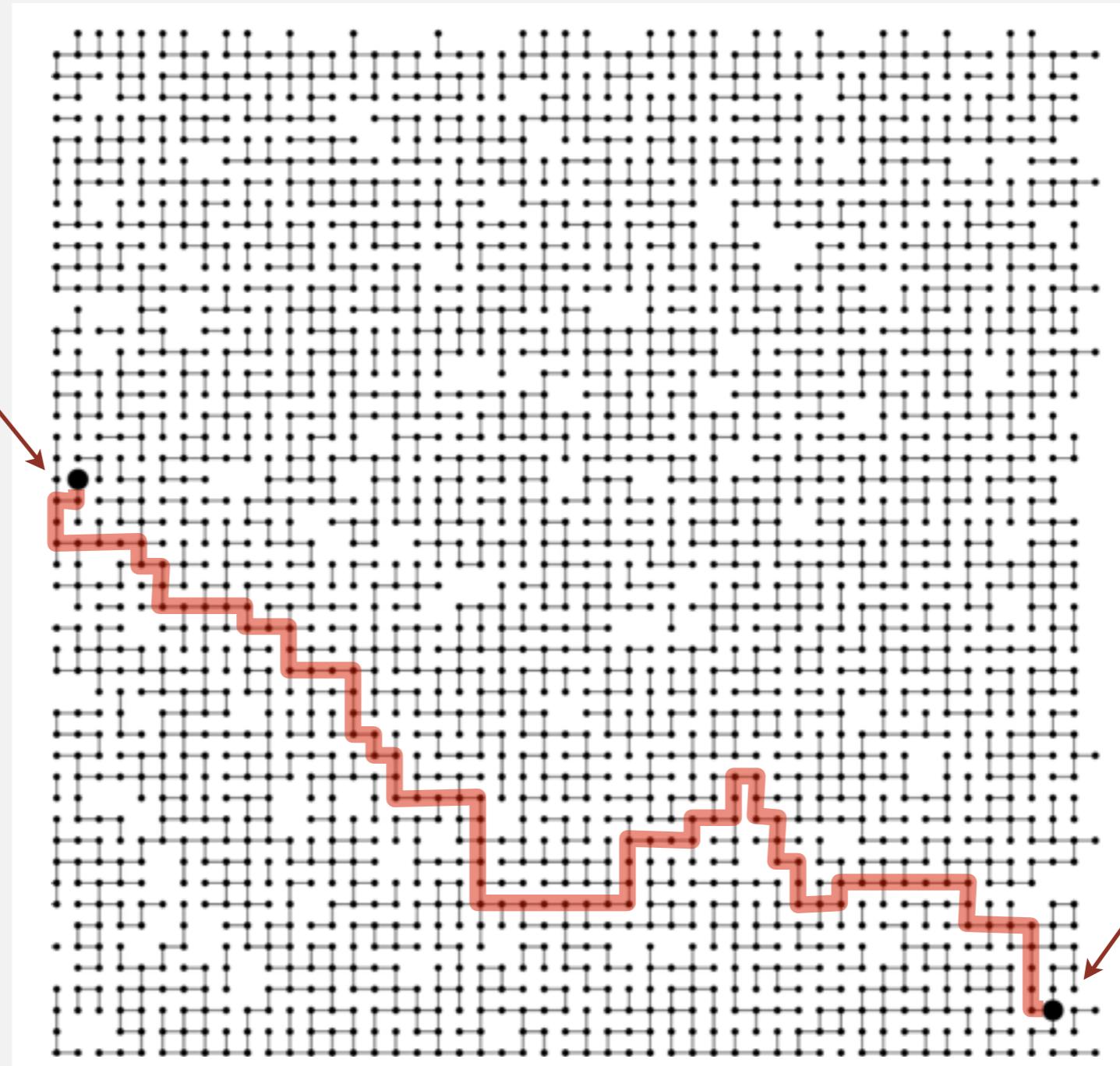
connect 1 and 0

are 0 and 7 connected? ✓



A larger connectivity example

Q. Is there a path connecting p and q ?



A. Yes.

Modeling the objects

Applications involve manipulating objects of all types.

- Pixels in a digital photo.
- Computers in a network.
- Friends in a social network.
- Transistors in a computer chip.
- Elements in a mathematical set.
- Variable names in a program.
- Metallic sites in a composite system.

When programming, convenient to name objects 0 to $N - 1$.

- Use integers as array index.
- Suppress details not relevant to union-find.



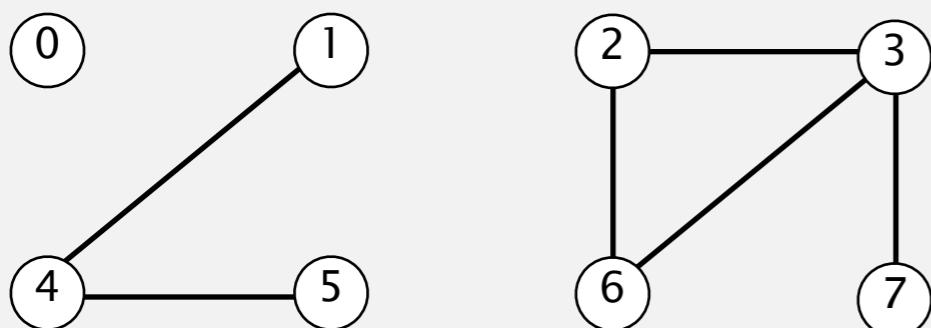
can use symbol table to translate from site names to integers: stay tuned (Chapter 3)

Modeling the connections

We assume "is connected to" is an equivalence relation:

- 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 connected to r ,
then p is connected to r .

Connected component. Maximal set of objects that are mutually connected.



{ 0 } { 1 4 5 } { 2 3 6 7 }

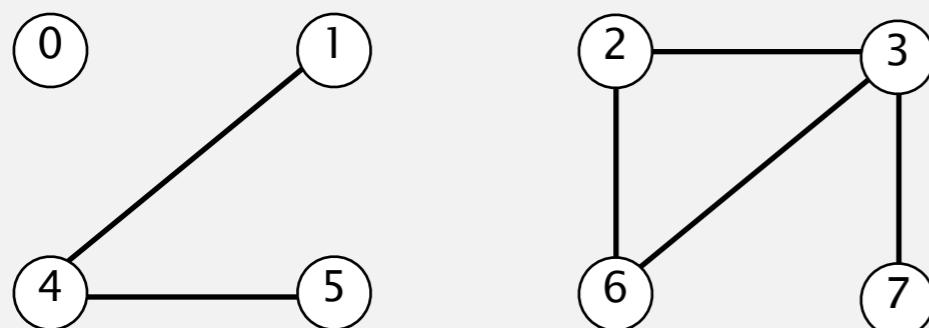
3 connected components

Implementing the operations

Find. In which component is object p ?

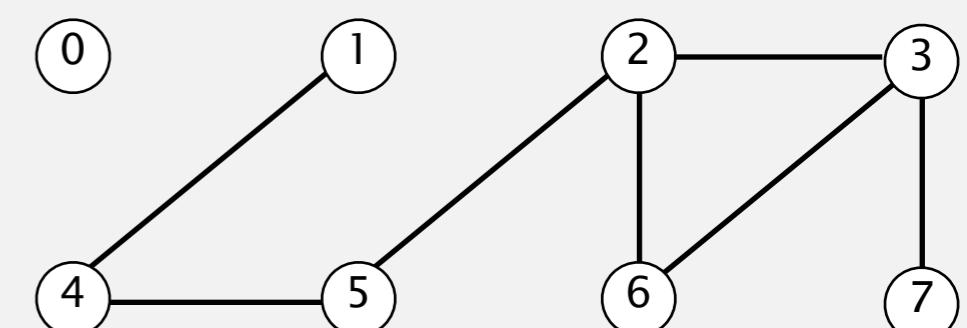
Connected. Are objects p and q in the same component?

Union. Replace components containing objects p and q with their union.



union(2, 5)

A large red arrow points from the initial graph state to the resulting graph state after the union operation.



{ 0 } { 1 4 5 } { 2 3 6 7 }

3 connected components

{ 0 } { 1 2 3 4 5 6 7 }

2 connected components

Union-find data type (API)

Goal. Design efficient data structure for union-find.

- Number of objects N can be huge.
- Number of operations M can be huge.
- Union and find operations may be intermixed.

```
public class UF
```

```
UF(int N)
```

*initialize union-find data structure
with N singleton objects (0 to $N - 1$)*

```
void union(int p, int q)
```

add connection between p and q

```
int find(int p)
```

component identifier for p (0 to $N - 1$)

```
boolean connected(int p, int q)
```

are p and q in the same component?

```
public boolean connected(int p, int q)
{ return find(p) == find(q); }
```

1-line implementation of connected()

Dynamic-connectivity client

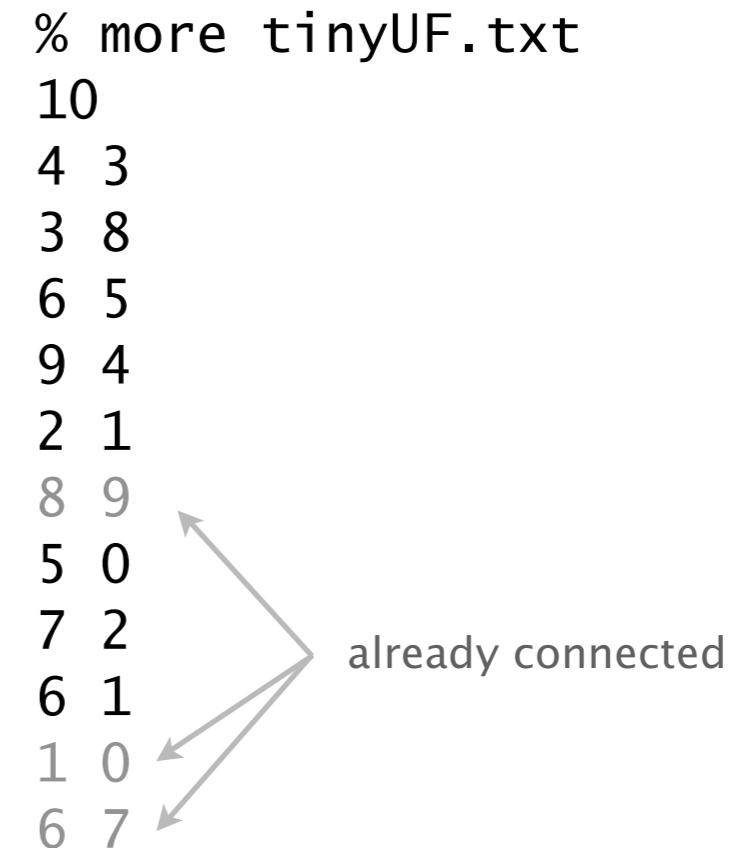
- Read in number of objects N from standard input.
- Repeat:
 - read in pair of integers from standard input
 - if they are not yet connected, connect them and print out pair

```
public static void main(String[] args)
{
    int N = StdIn.readInt();
    UF uf = new UF(N);
    while (!StdIn.isEmpty())
    {
        int p = StdIn.readInt();
        int q = StdIn.readInt();
        if (!uf.connected(p, q))
        {
            uf.union(p, q);
            StdOut.println(p + " " + q);
        }
    }
}
```

% more tinyUF.txt

10	
4	3
3	8
6	5
9	4
2	1
8	9
5	0
7	2
6	1
1	0
6	7

already connected



Algorithms

ROBERT SEDGEWICK | KEVIN WAYNE

<http://algs4.cs.princeton.edu>

1.5 UNION-FIND

- ▶ *dynamic connectivity*
- ▶ ***quick find***
- ▶ *quick union*
- ▶ *improvements*
- ▶ *applications*

Quick-find [eager approach]

Data structure.

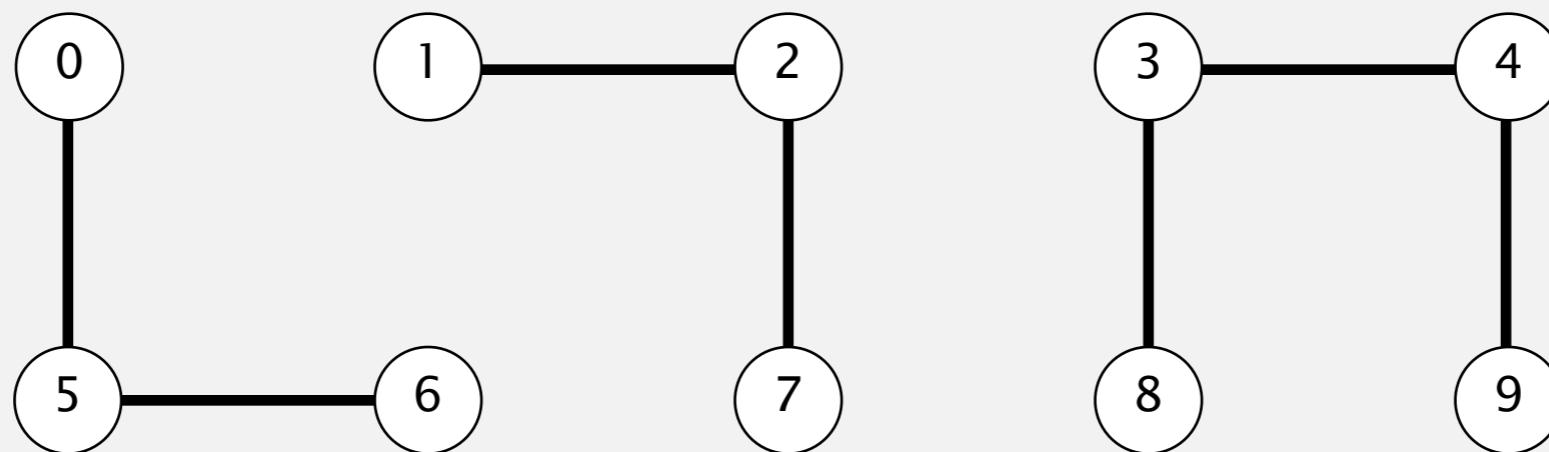
- Integer array $\text{id}[]$ of length N .
- Interpretation: $\text{id}[p]$ is the id of the component containing p .

if and only if

↙

	0	1	2	3	4	5	6	7	8	9
$\text{id}[]$	0	1	1	8	8	0	0	1	8	8

0, 5 and 6 are connected
1, 2, and 7 are connected
3, 4, 8, and 9 are connected



Quick-find [eager approach]

Data structure.

- Integer array $\text{id}[]$ of length N .
- Interpretation: $\text{id}[p]$ is the id of the component containing p .

	0	1	2	3	4	5	6	7	8	9
$\text{id}[]$	0	1	1	8	8	0	0	1	8	8

Find. What is the id of p ?

$\text{id}[6] = 0; \text{id}[1] = 1$
6 and 1 are not connected

Connected. Do p and q have the same id?

Union. To merge components containing p and q , change all entries whose id equals $\text{id}[p]$ to $\text{id}[q]$.

	0	1	2	3	4	5	6	7	8	9
$\text{id}[]$	1	1	1	8	8	1	1	1	8	8

problem: many values can change

after union of 6 and 1

Quick-find demo



0

1

2

3

4

5

6

7

8

9

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

Quick-find demo

union(4, 3)

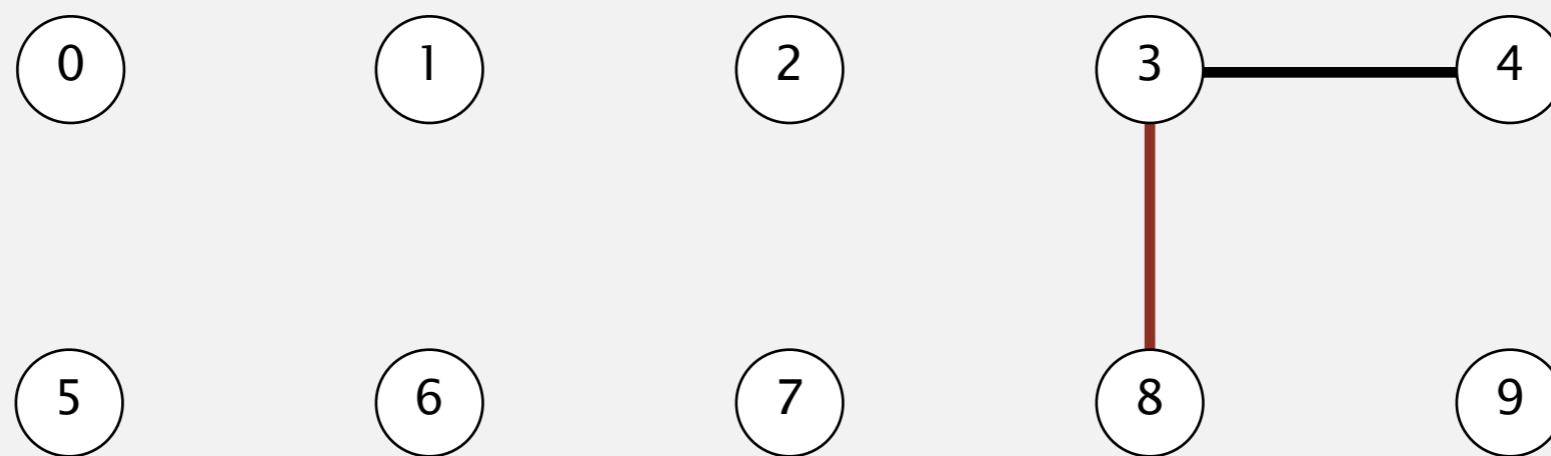


id[]	0	1	2	3	4	5	6	7	8	9
	0	1	2	3	3	5	6	7	8	9

Two red arrows point to the value 3 in the array, indicating the union operation is being performed between index 4 and index 3.

Quick-find demo

union(3, 8)

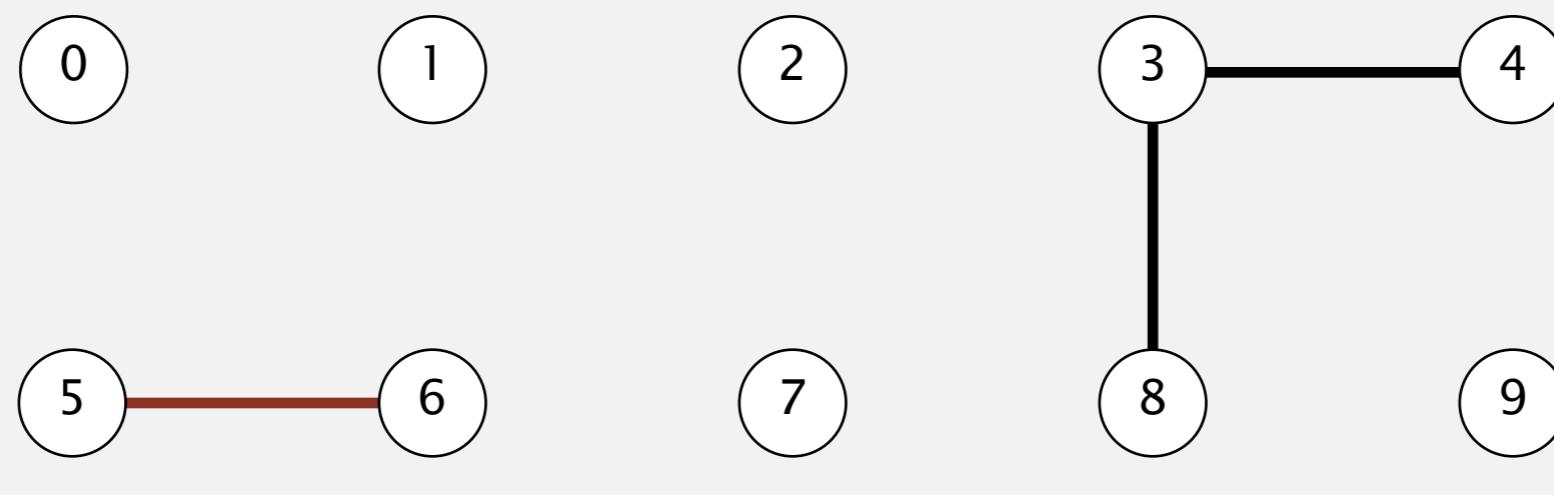


	0	1	2	3	4	5	6	7	8	9
id[]	0	1	2	8	8	5	6	7	8	9

↑ ↑

Quick-find demo

union(6, 5)

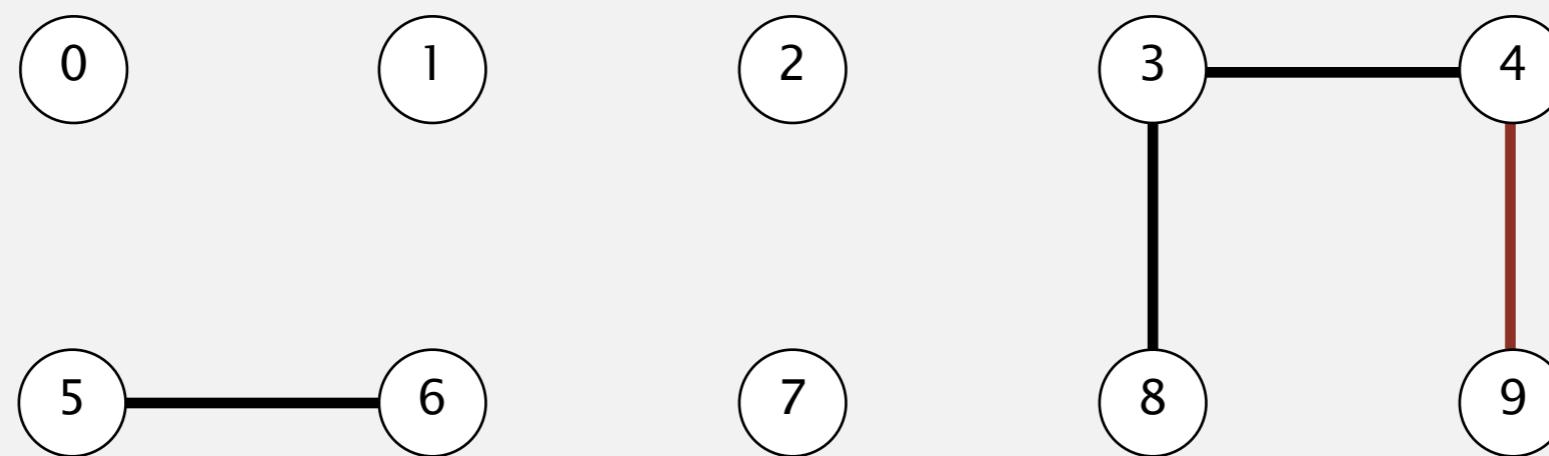


	0	1	2	3	4	5	6	7	8	9
id[]	0	1	2	8	8	5	5	7	8	9

↑ ↑

Quick-find demo

union(9, 4)

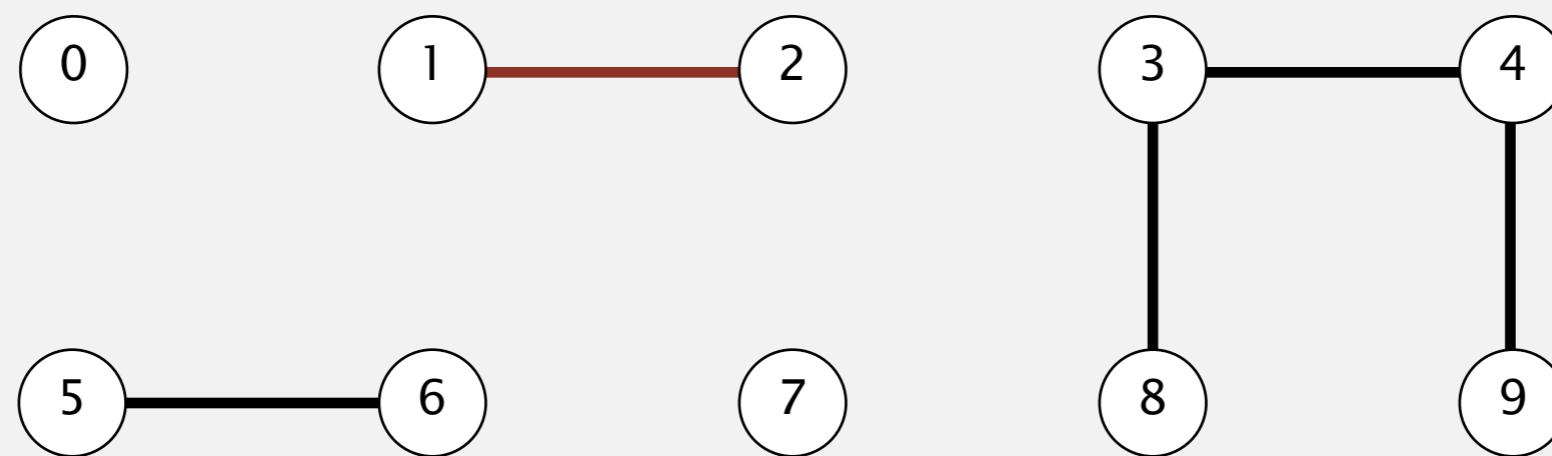


	0	1	2	3	4	5	6	7	8	9
id[]	0	1	2	8	8	5	5	7	8	8

Arrows point from the 4th and 9th indices of the id[] array to the values 8 and 8 respectively.

Quick-find demo

union(2, 1)

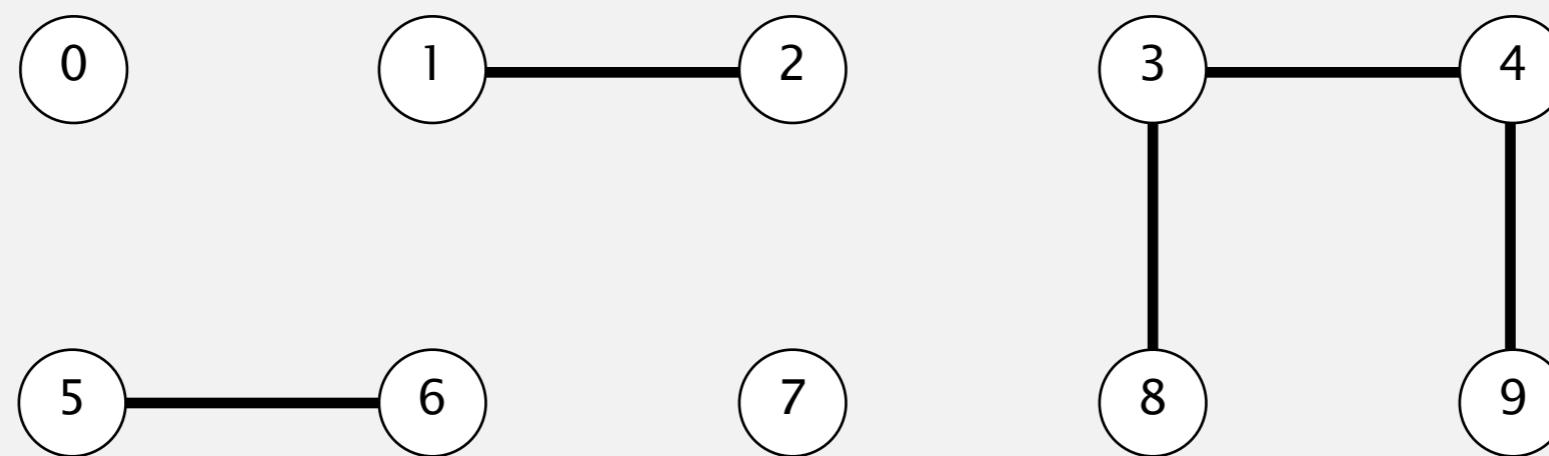


	0	1	2	3	4	5	6	7	8	9
id[]	0	1	1	8	8	5	5	7	8	8

↑ ↑

Quick-find demo

connected(8, 9)



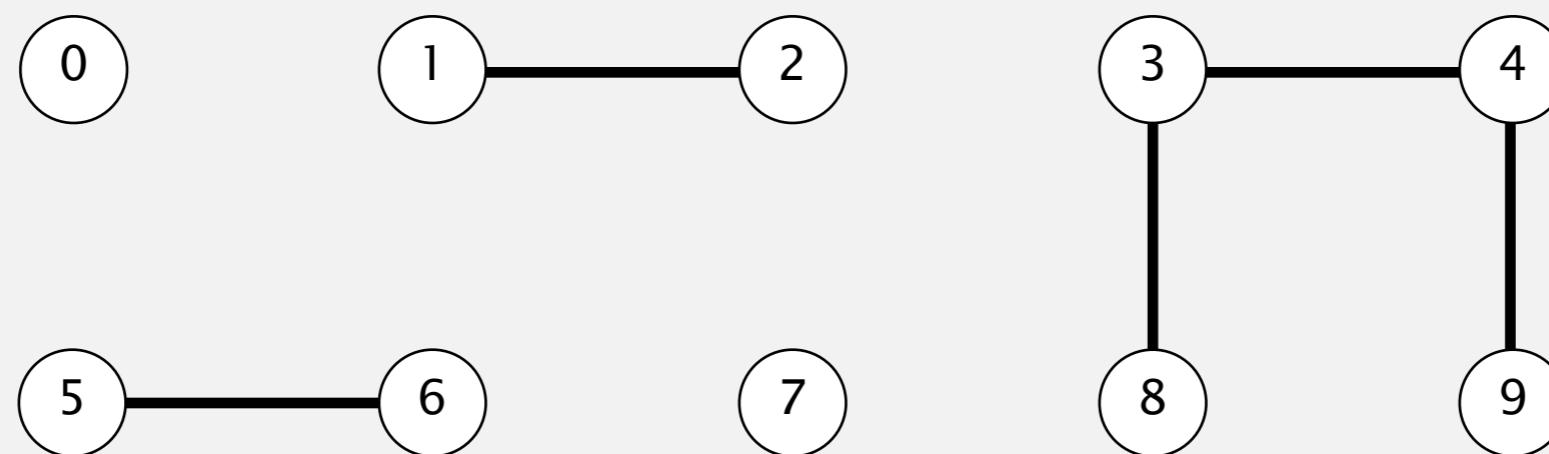
	0	1	2	3	4	5	6	7	8	9
id[]	0	1	1	8	8	5	5	7	8	8

↑ ↑

already connected

Quick-find demo

connected(5, 0)

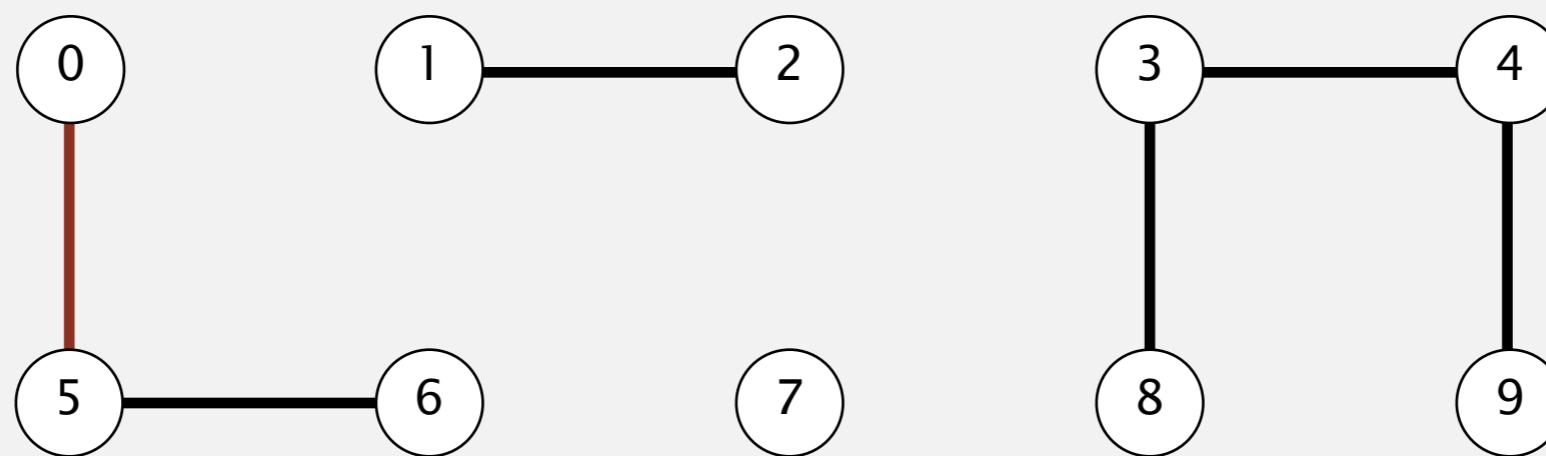


	0	1	2	3	4	5	6	7	8	9
id[]	0	1	1	8	8	5	5	7	8	8
	↑			↑						

not connected

Quick-find demo

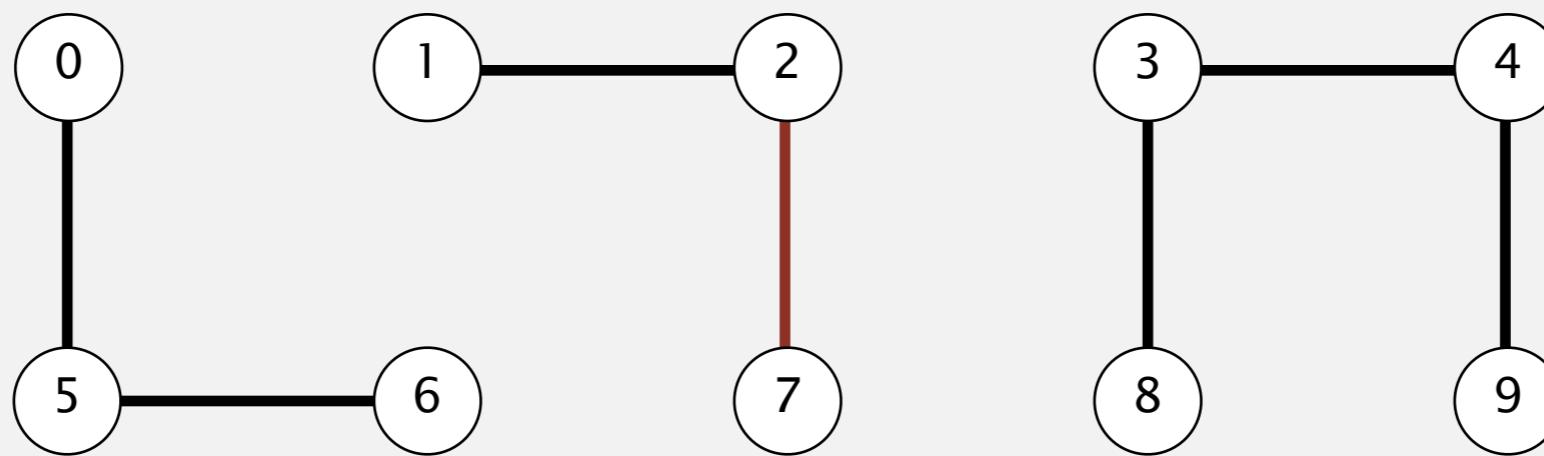
union(5, 0)



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

Quick-find demo

union(7, 2)

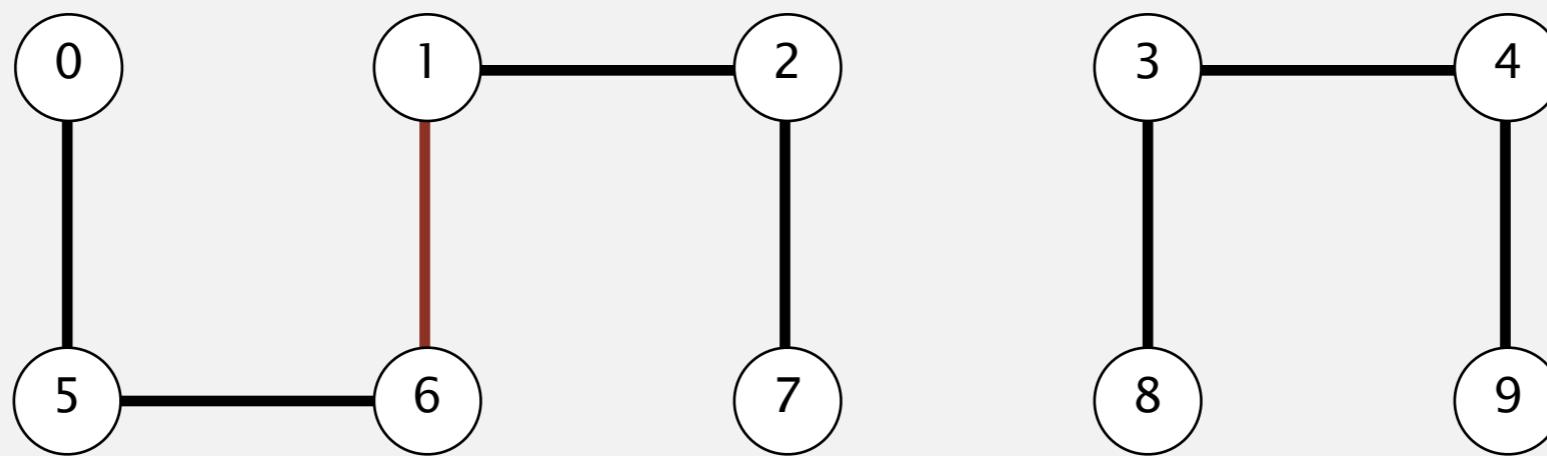


	0	1	2	3	4	5	6	7	8	9
id[]	0	1	1	8	8	0	0	1	8	8

Arrows point to the '1' at index 7 and the '1' at index 7.

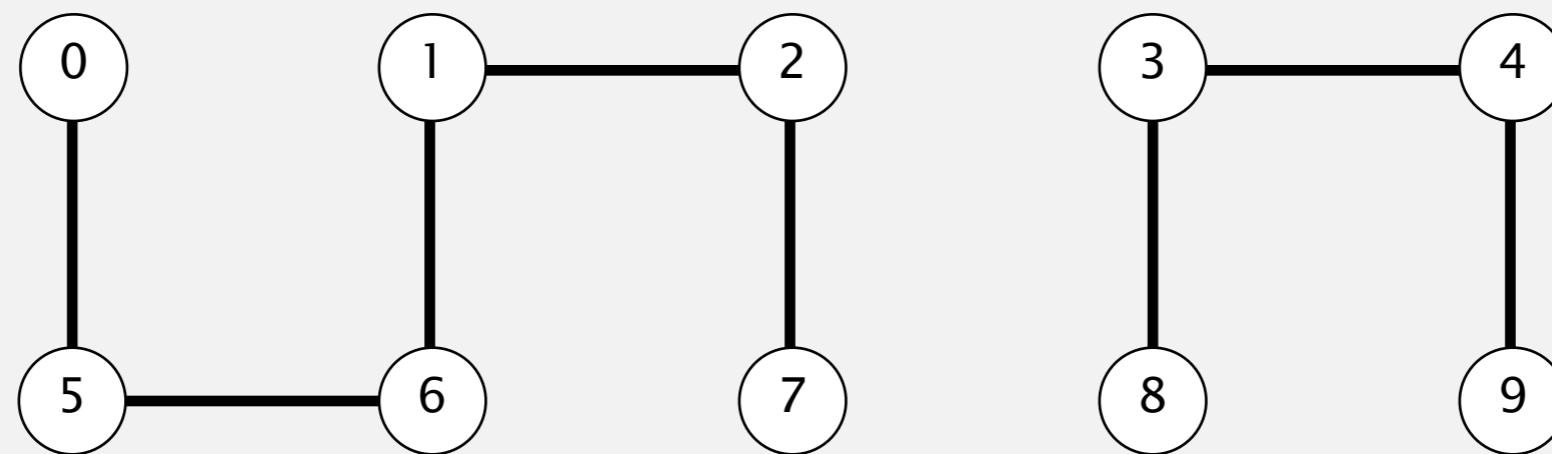
Quick-find demo

union(6, 1)



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

Quick-find demo



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

Quick-find: Java implementation

```
public class QuickFindUF
{
    private int[] id;
```

```
    public QuickFindUF(int N)
    {
```

```
}
```

```
    public boolean find(int p)
```

```
    public void union(int p, int q)
    {
```

```
}
```



set id of each object to itself
(N array accesses)



return the id of p
(1 array access)



change all entries with $\text{id}[p]$ to $\text{id}[q]$
(at most $2N + 2$ array accesses)

Quick-find: Java implementation

```
public class QuickFindUF
{
    private int[] id;
```

```
public QuickFindUF(int N)
{
```

```
    id = new int[N];
    for (int i = 0; i < N; i++)
        id[i] = i;
```

```
}
```

set id of each object to itself
(N array accesses)

```
public boolean find(int p)
{    return id[p]; }
```

return the id of p
(1 array access)

```
public void union(int p, int q)
{
```

```
    int pid = id[p];
    int qid = id[q];
    for (int i = 0; i < id.length; i++)
        if (id[i] == pid) id[i] = qid;
```

```
}
```

```
}
```

change all entries with $\text{id}[p]$ to $\text{id}[q]$
(at most $2N + 2$ array accesses)

Quick-find is too slow

Cost model. Number of array accesses (for read or write).

algorithm	initialize	union	find	connected
quick-find	N	N	1	1

order of growth of number of array accesses

Union is too expensive. It takes N^2 array accesses to process a sequence of N union operations on N objects.

quadratic

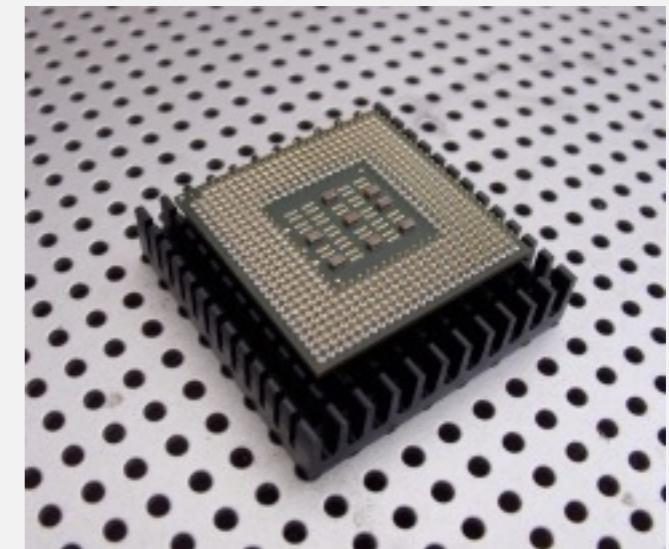


Quadratic algorithms do not scale

Rough standard (for now).

- 10^9 operations per second.
- 10^9 words of main memory.
- Touch all parts of memory in approx. 1 second.

a truism (roughly)
since 1950!

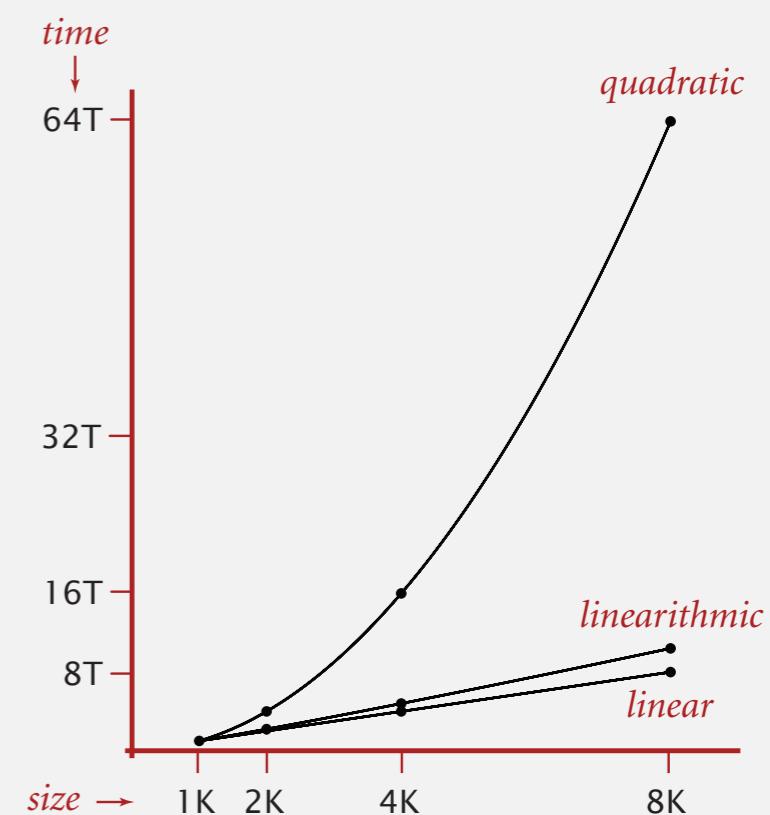


Ex. Huge problem for quick-find.

- 10^9 union commands on 10^9 objects.
- Quick-find takes more than 10^{18} operations.
- 30+ years of computer time!

Quadratic algorithms don't scale with technology.

- New computer may be 10x as fast.
- But, has 10x as much memory ⇒ want to solve a problem that is 10x as big.
- With quadratic algorithm, takes 10x as long!



Algorithms

ROBERT SEDGEWICK | KEVIN WAYNE

<http://algs4.cs.princeton.edu>

1.5 UNION-FIND

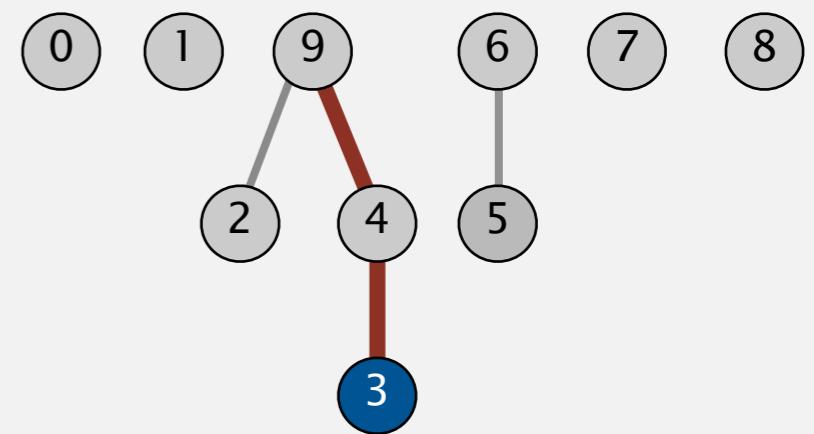
- ▶ *dynamic connectivity*
- ▶ *quick find*
- ▶ *quick union*
- ▶ *improvements*
- ▶ *applications*

Quick-union [lazy approach]

Data structure.

- Integer array $\text{id}[]$ of length N .
- Interpretation: $\text{id}[i]$ is parent of i . keep going until it doesn't change
(algorithm ensures no cycles)
- Root of i is $\text{id}[\text{id}[\text{id}[\dots \text{id}[i]\dots]]]$.

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



parent of 3 is 4

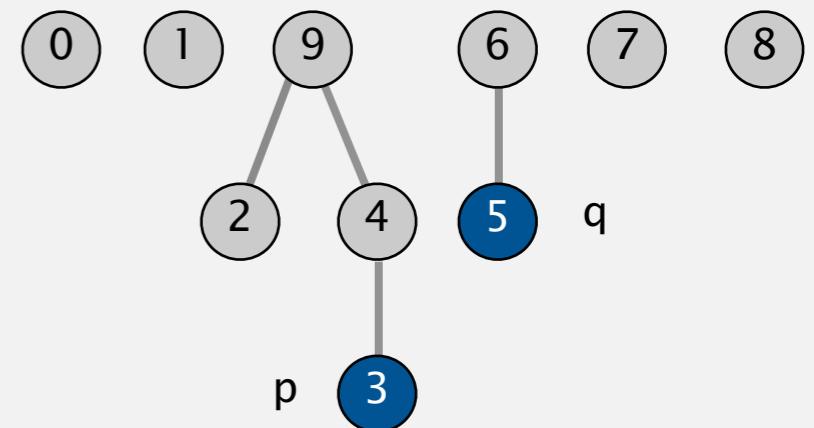
root of 3 is 9

Quick-union [lazy approach]

Data structure.

- Integer array $\text{id}[]$ of length N .
- Interpretation: $\text{id}[i]$ is parent of i .
- Root of i is $\text{id}[\text{id}[\text{id}[\dots \text{id}[i]\dots]]]$.

0	1	2	3	4	5	6	7	8	9	
$\text{id}[]$	0	1	9	4	9	6	6	7	8	9



Find. What is the root of p ?

Connected. Do p and q have the same root?

root of 3 is 9

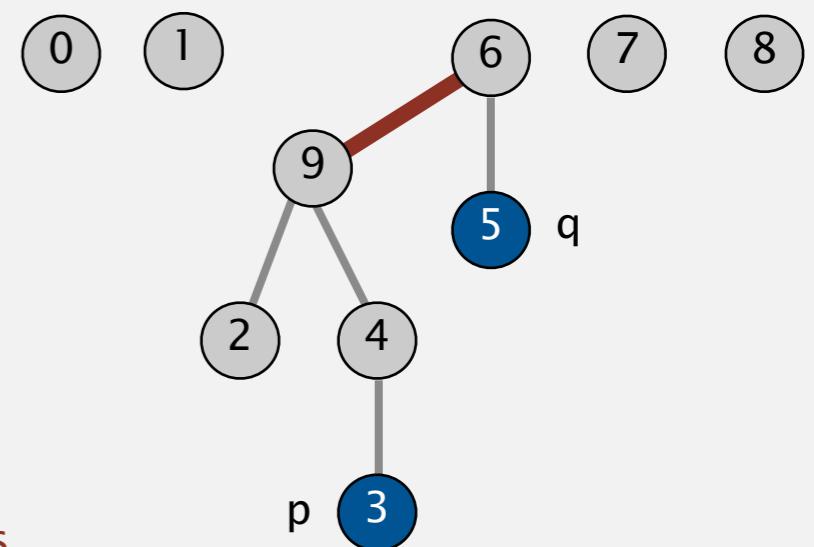
root of 5 is 6

3 and 5 are not connected

Union. To merge components containing p and q , set the id of p 's root to the id of q 's root.

0	1	2	3	4	5	6	7	8	9	
$\text{id}[]$	0	1	9	4	9	6	6	7	8	6

only one value changes



Quick-union demo



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

Quick-union demo

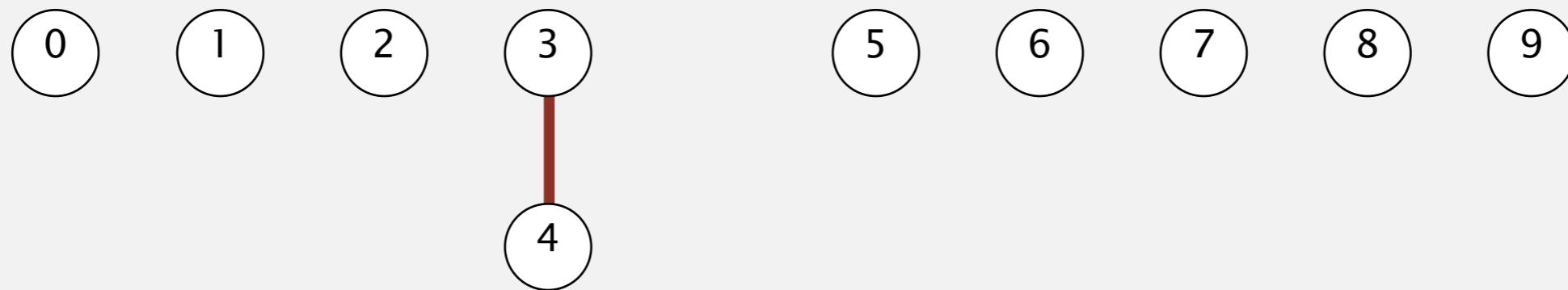
union(4, 3)



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

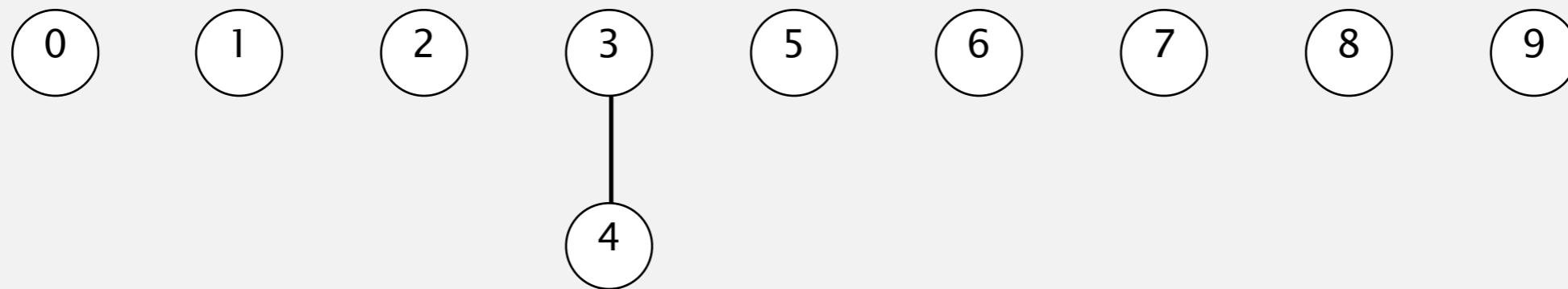
Quick-union demo

union(4, 3)



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

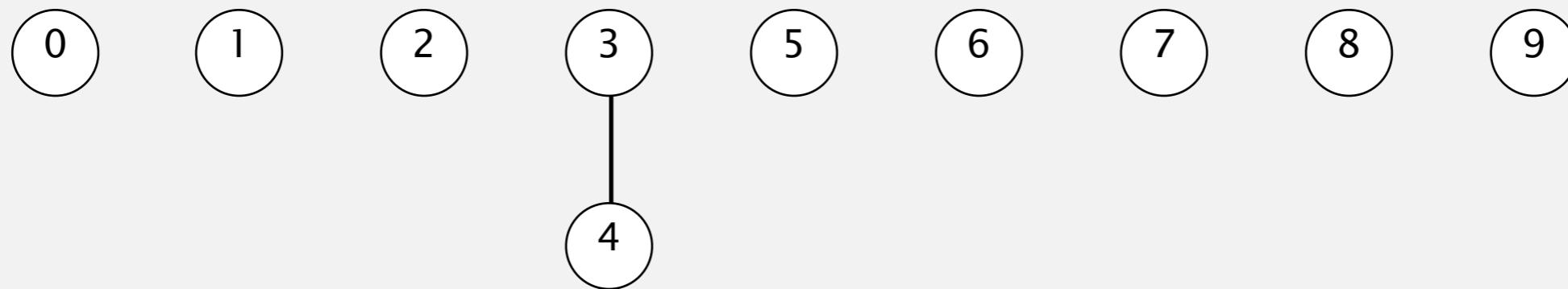
Quick-union demo



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

Quick-union demo

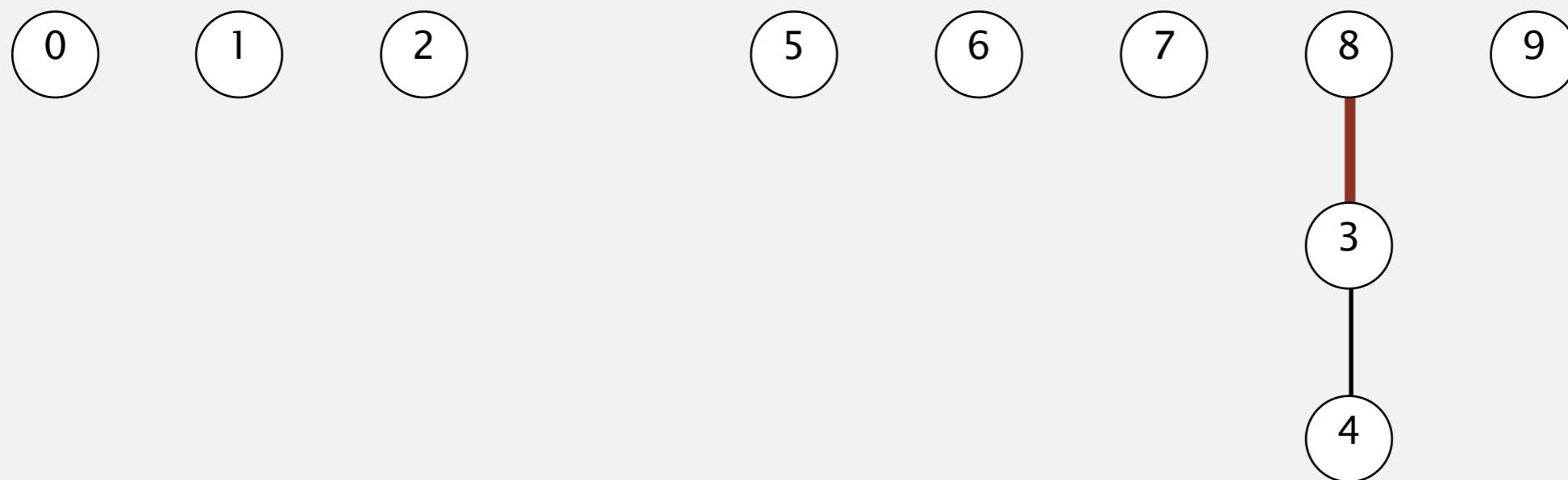
union(3, 8)



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

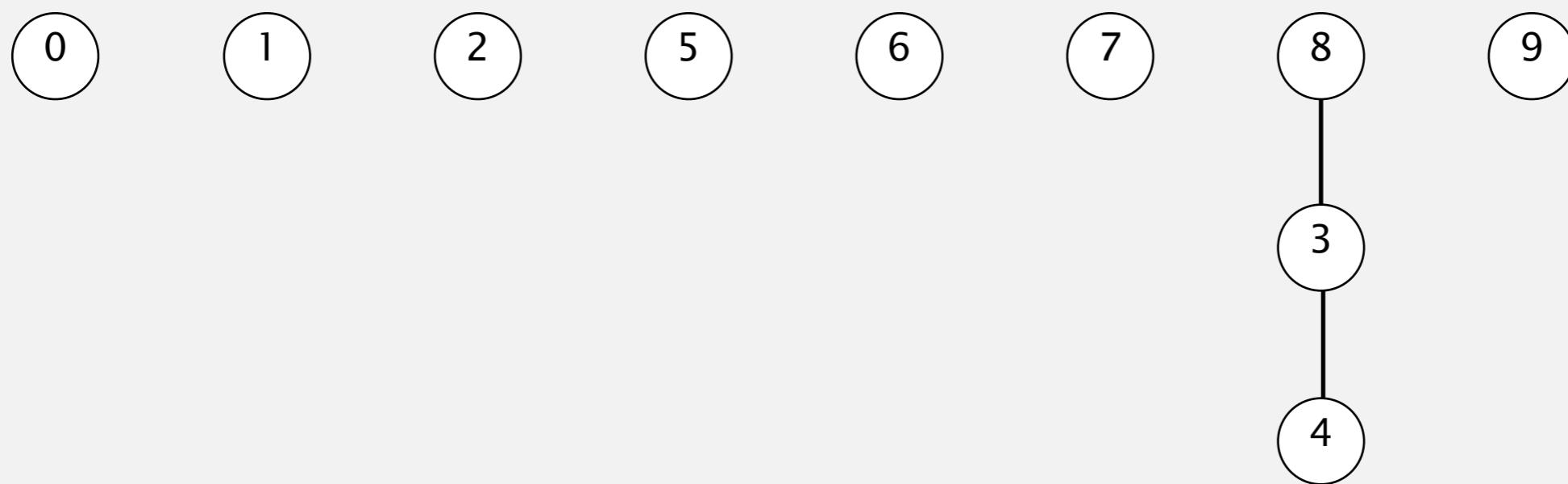
Quick-union demo

union(3, 8)



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

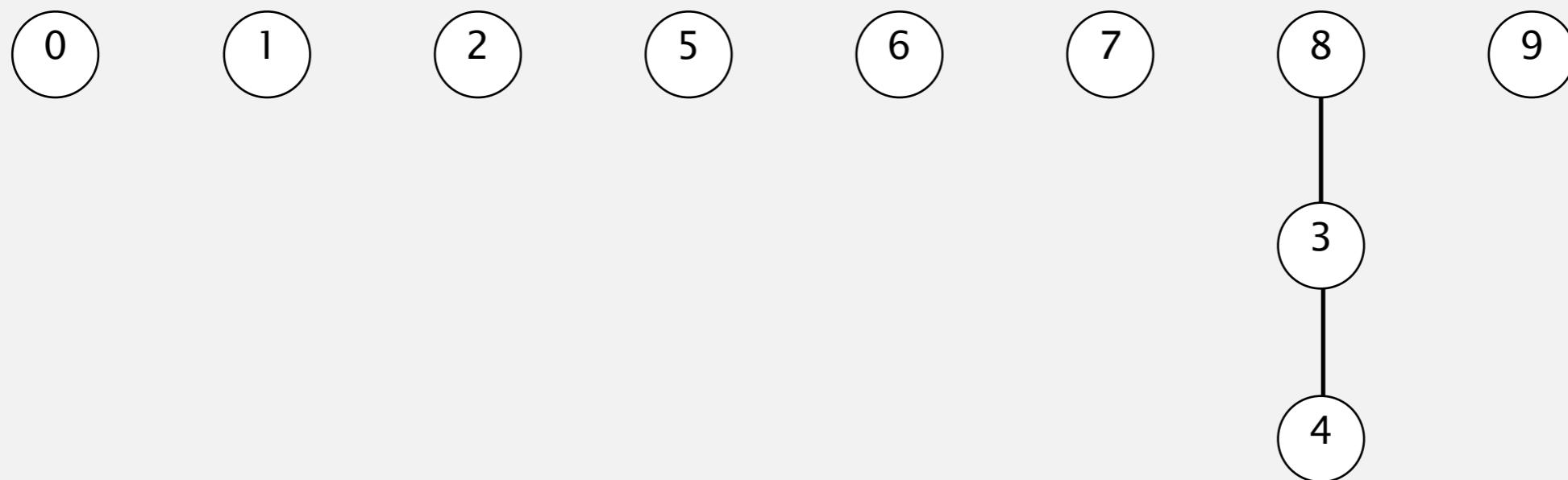
Quick-union demo



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

Quick-union demo

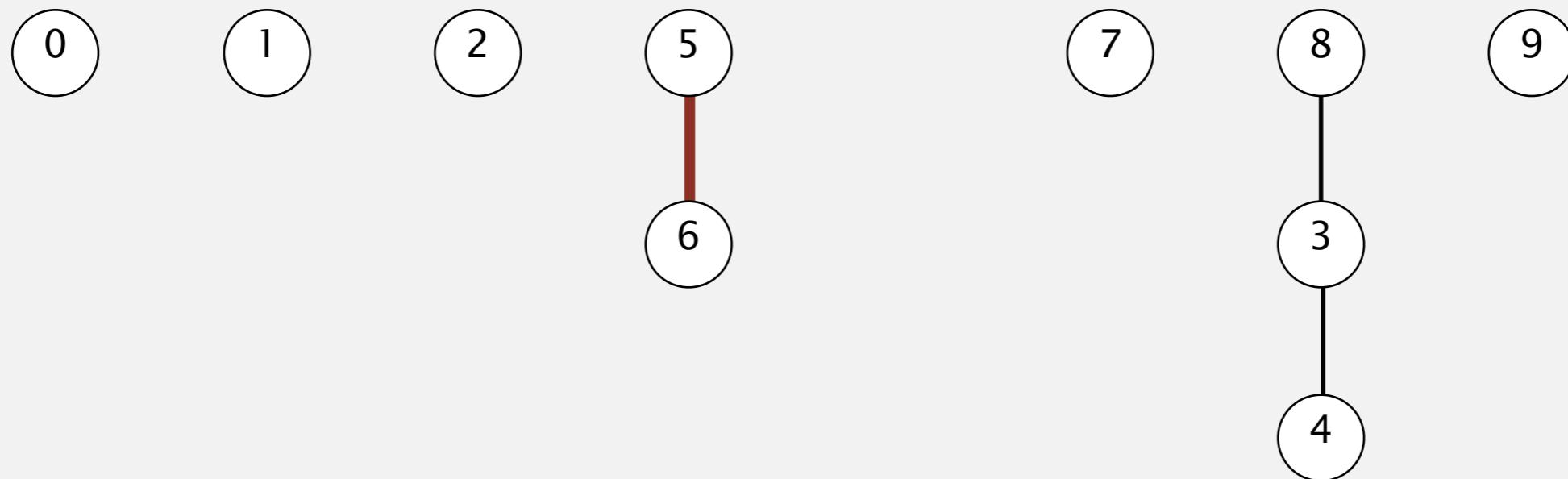
union(6, 5)



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

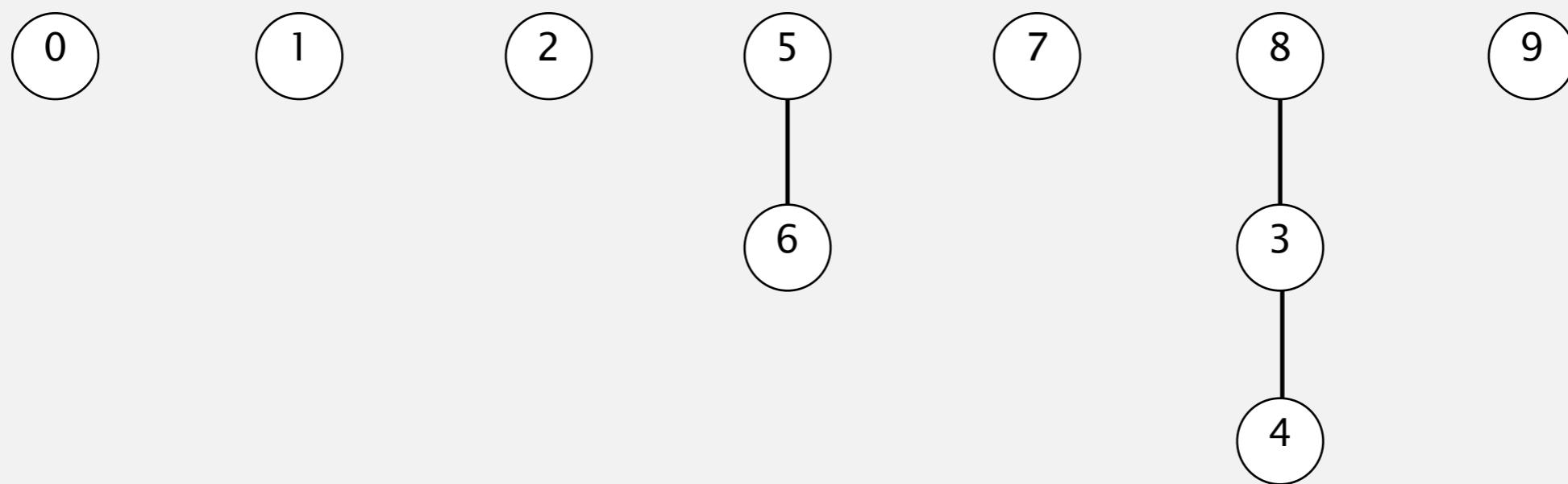
Quick-union demo

union(6, 5)



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

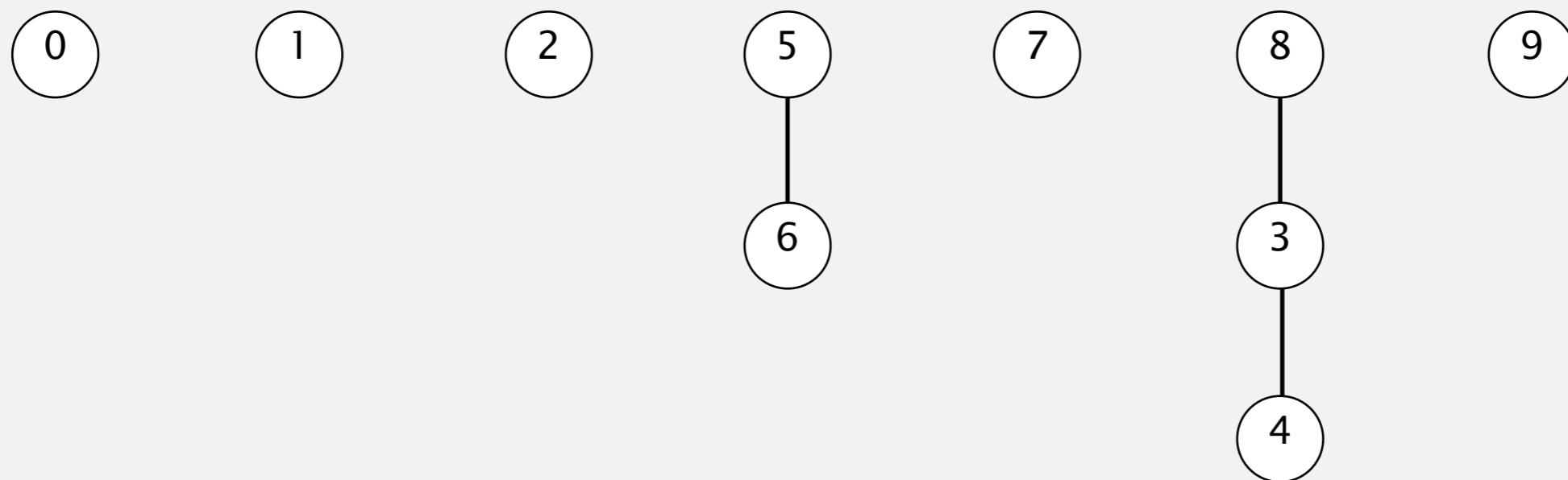
Quick-union demo



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

Quick-union demo

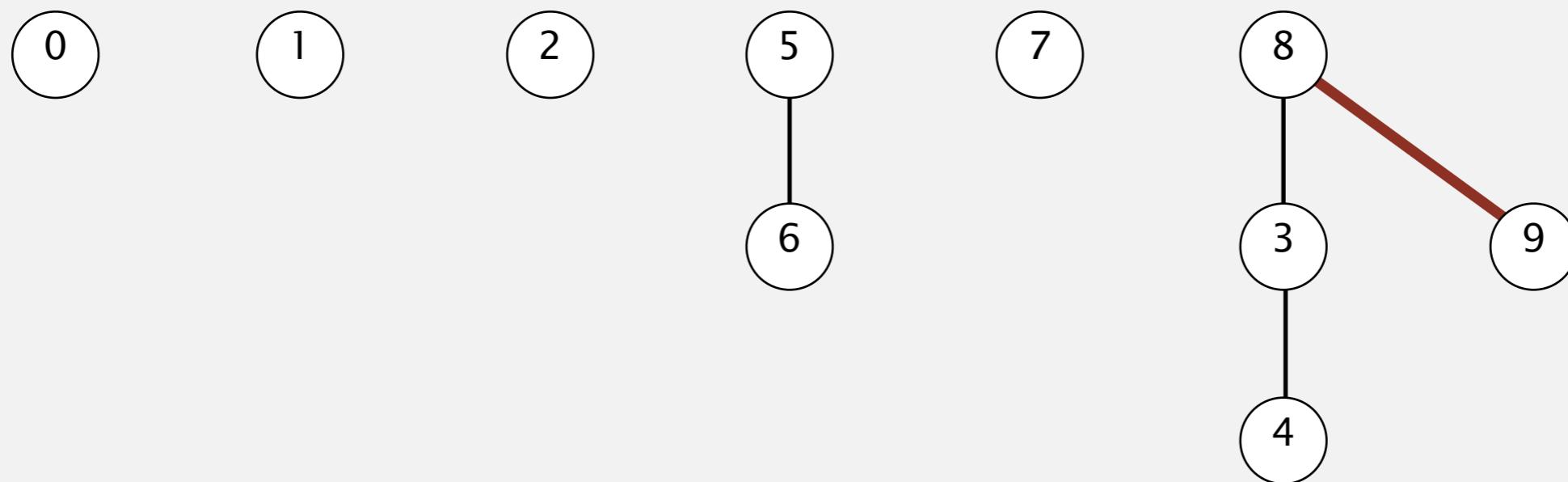
union(9, 4)



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

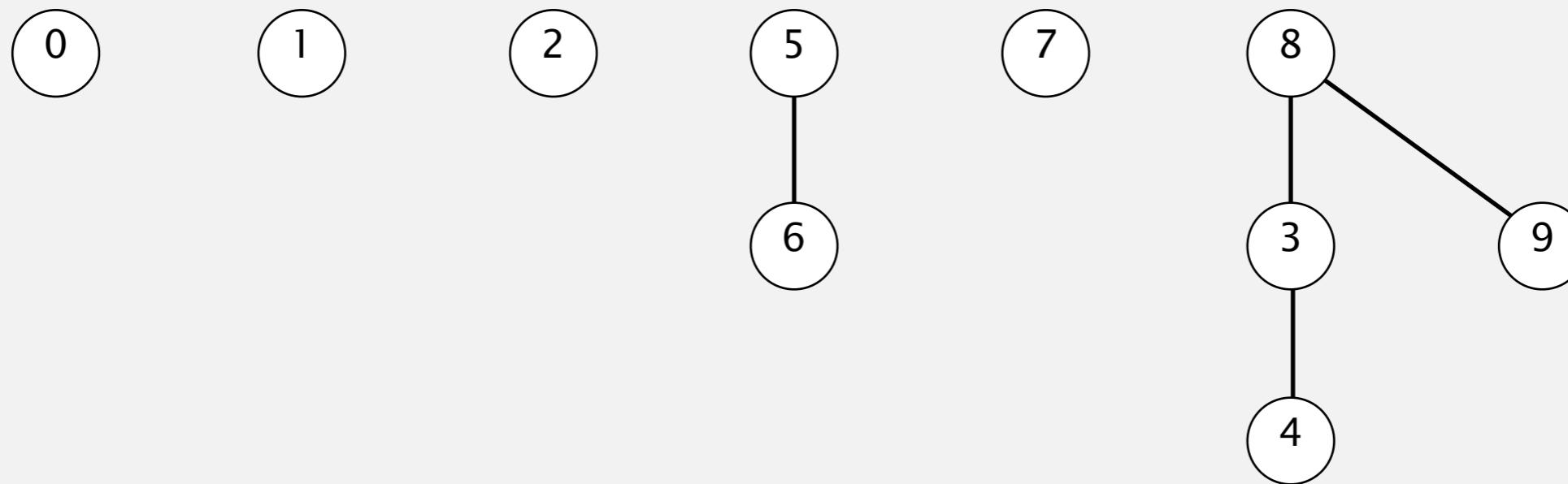
Quick-union demo

union(9, 4)



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

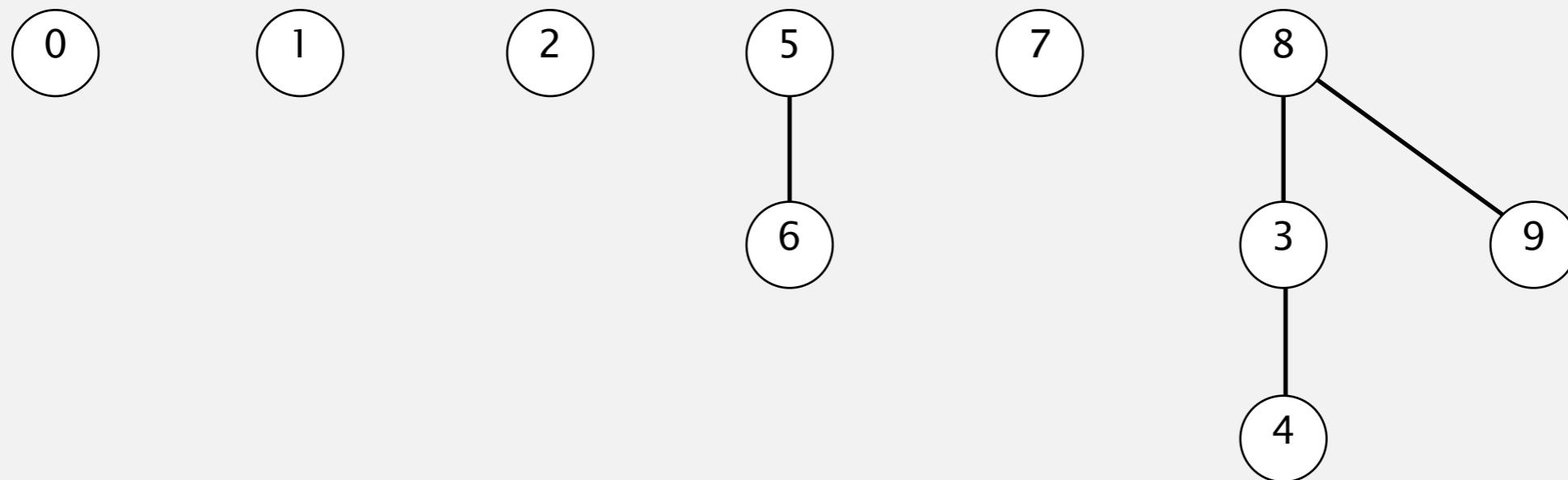
Quick-union demo



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

Quick-union demo

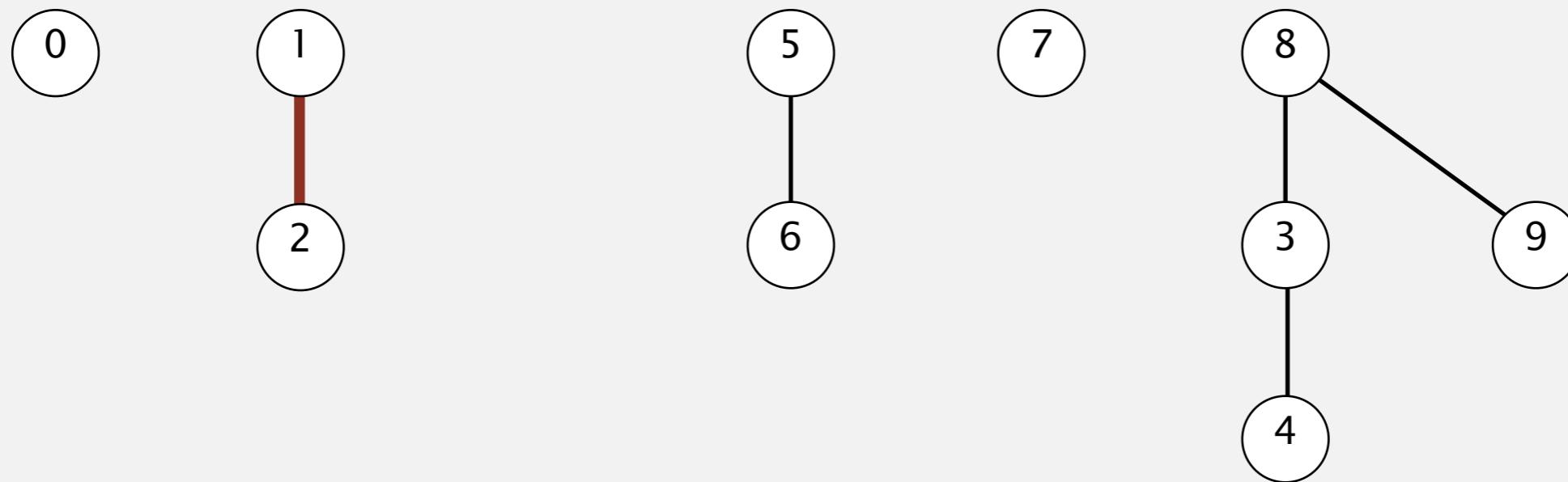
union(2, 1)



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

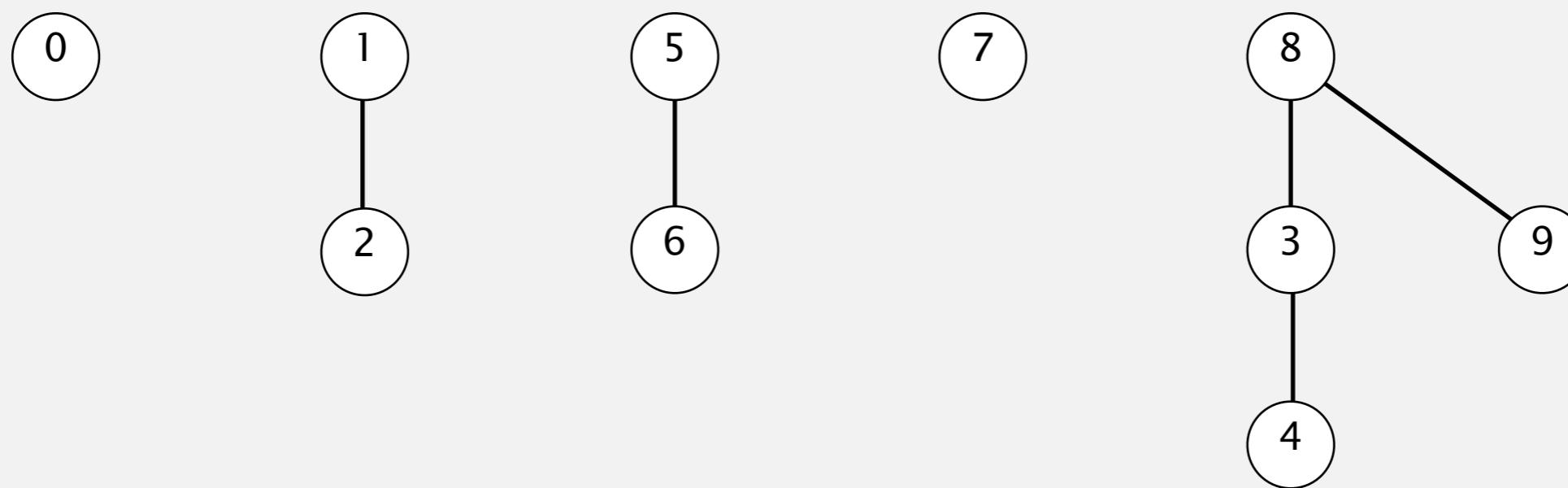
Quick-union demo

union(2, 1)



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

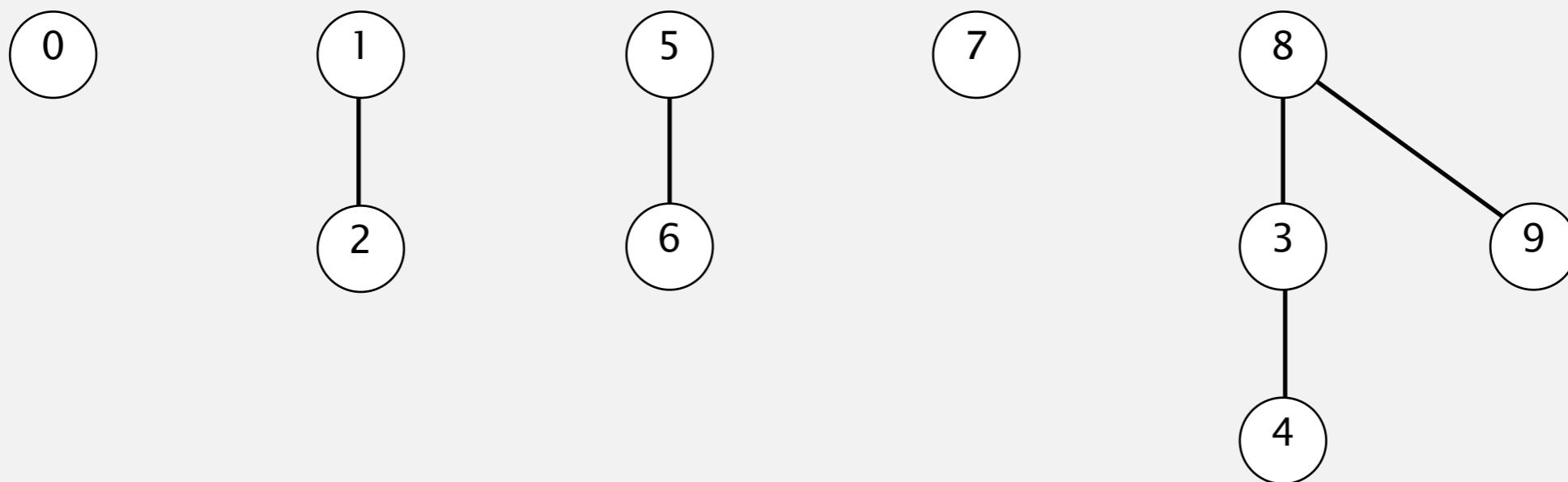
Quick-union demo



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

Quick-union demo

connected(8, 9) ✓

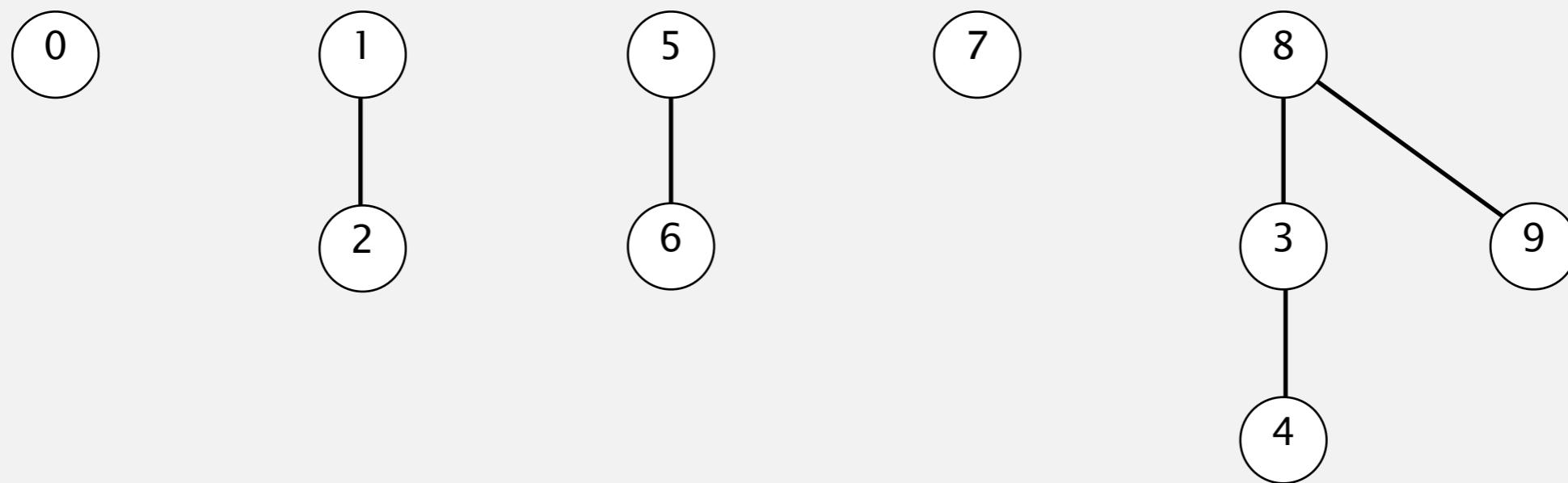


0	1	2	3	4	5	6	7	8	9
0	1	1	8	3	5	5	7	8	8

Quick-union demo

connected(5, 4)

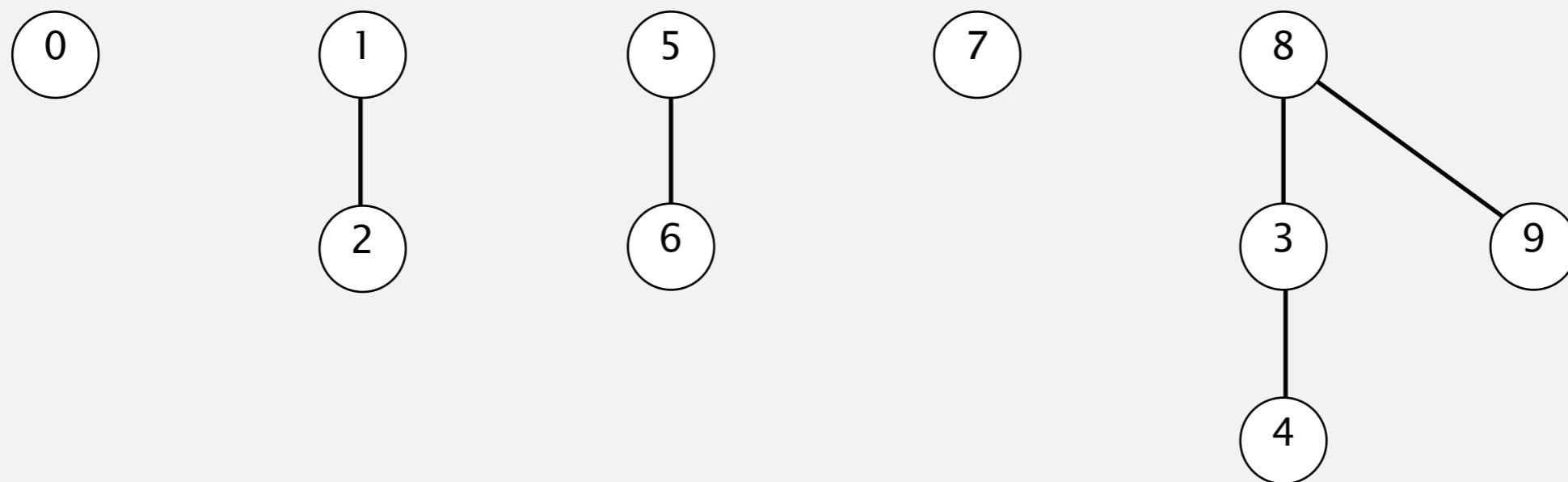
X



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

Quick-union demo

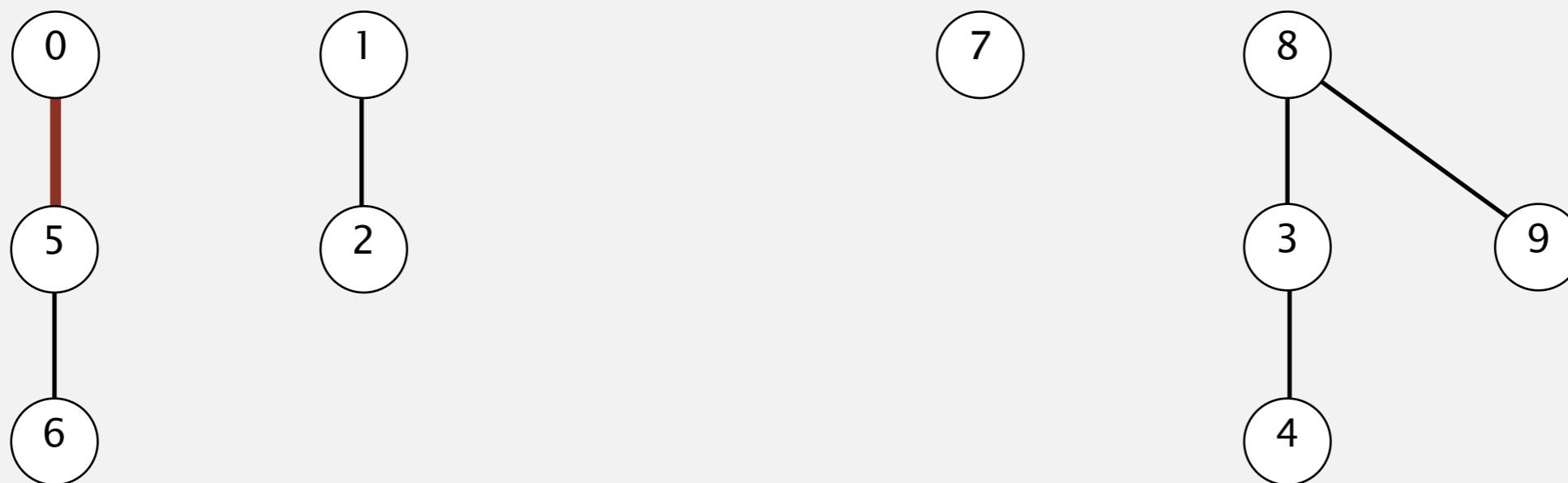
union(5, 0)



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

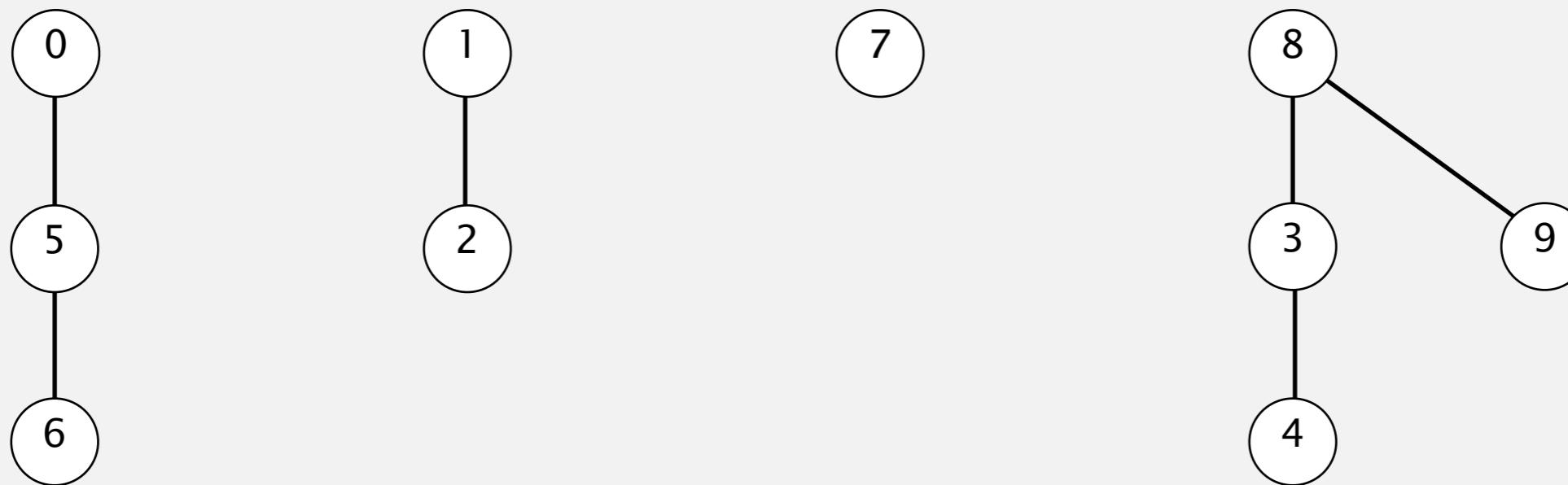
Quick-union demo

union(5, 0)



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

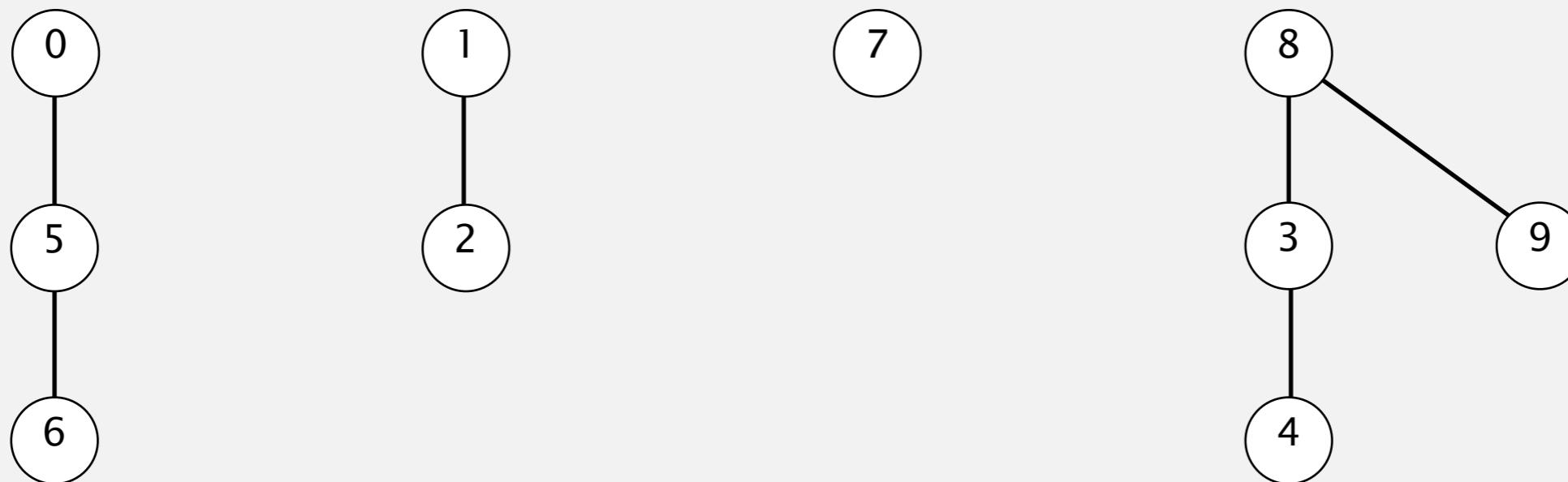
Quick-union demo



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

Quick-union demo

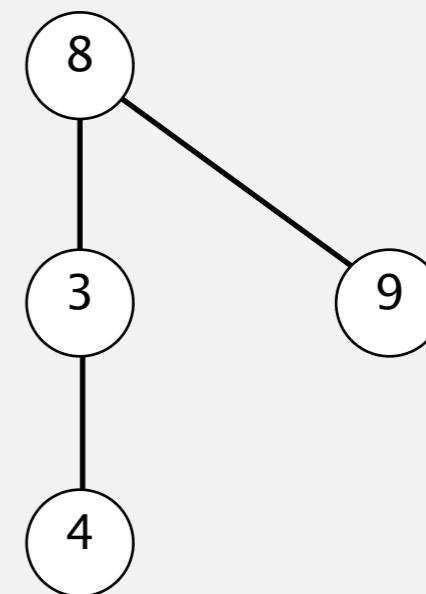
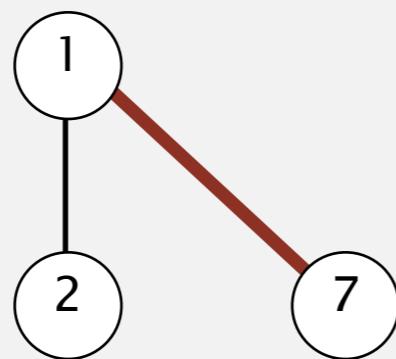
union(7, 2)



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

Quick-union demo

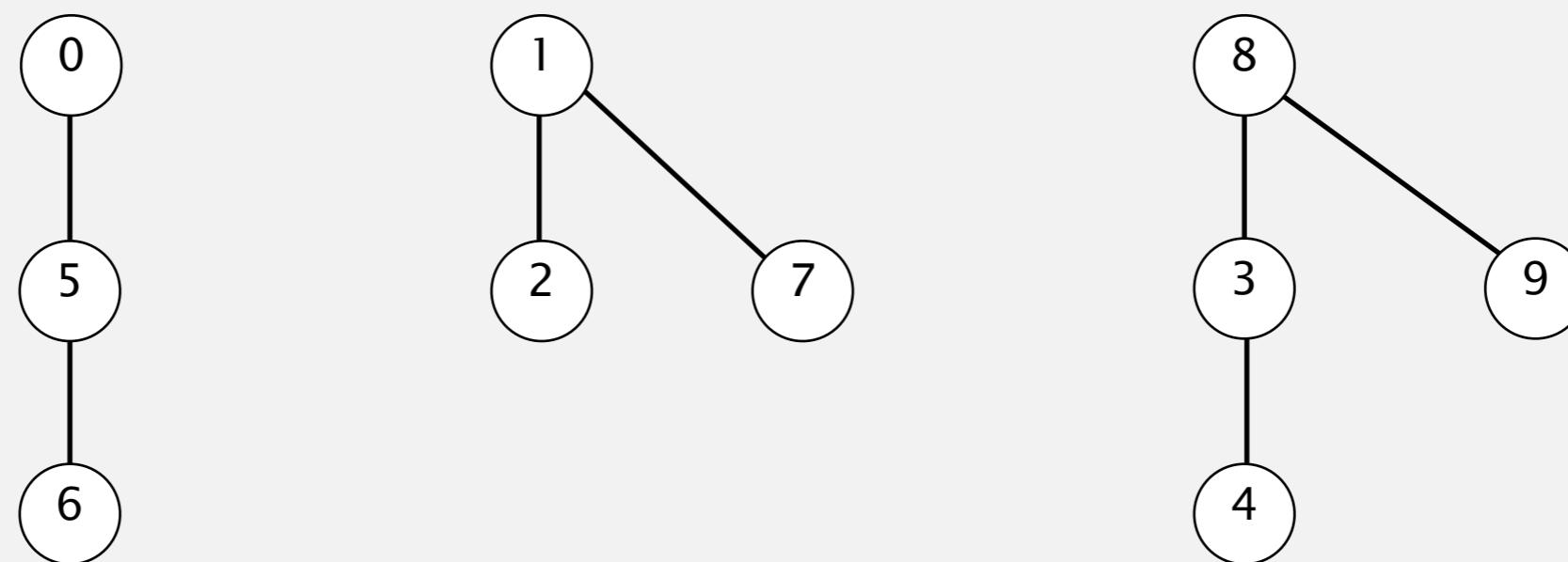
union(7, 2)



0 1 2 3 4 5 6 7 8 9

id[]	0	1	1	8	3	0	5	1	8	8
	0	1	1	8	3	0	5	1	8	8

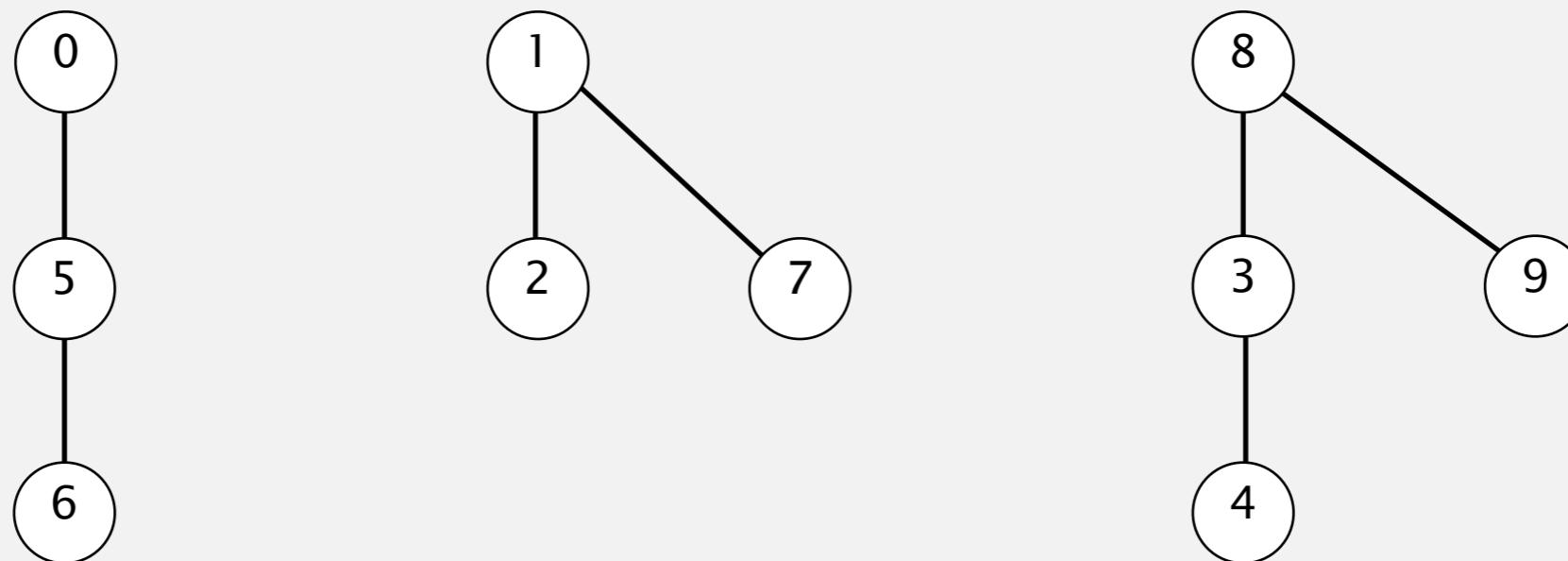
Quick-union demo



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

Quick-union demo

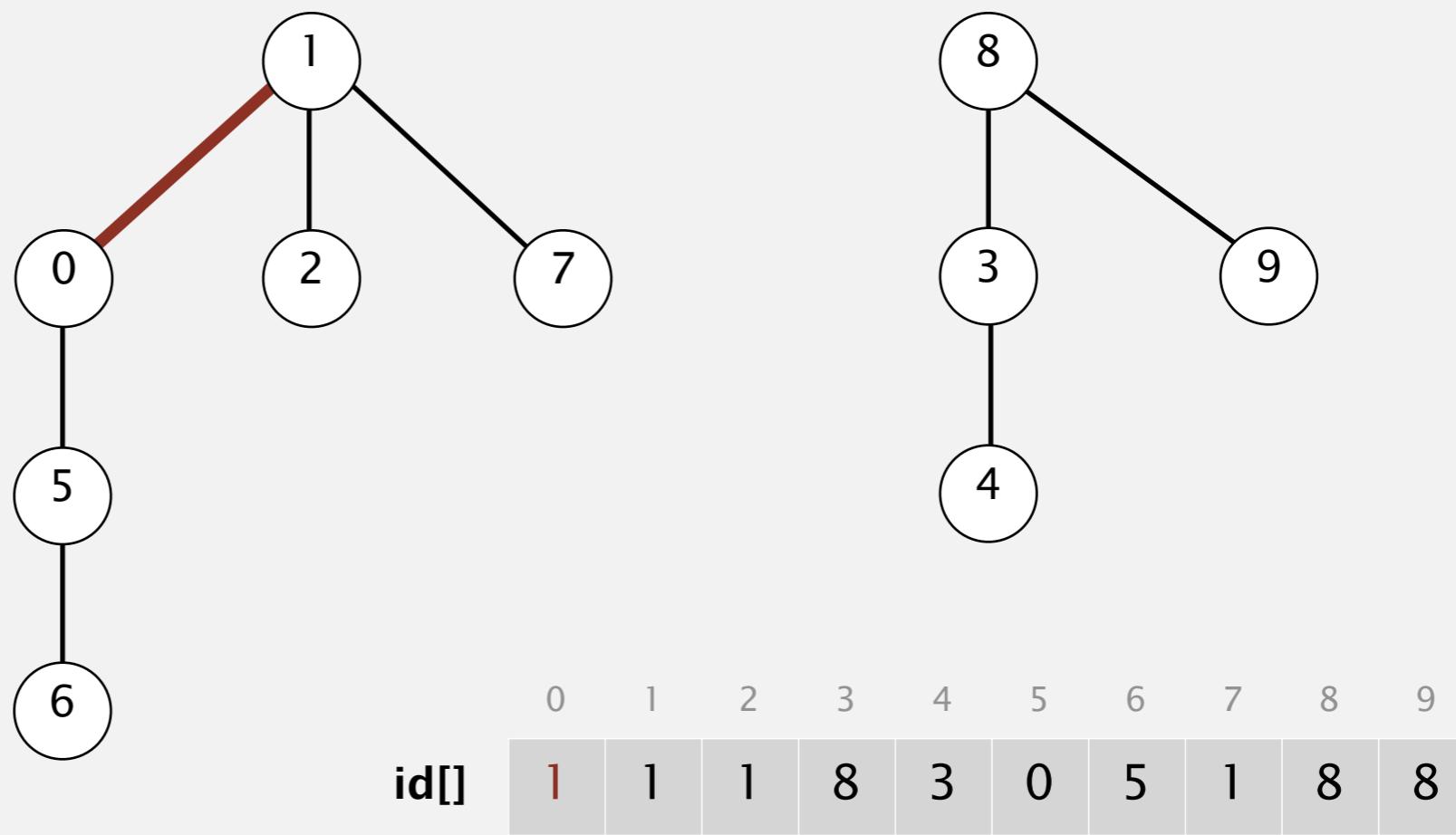
union(6, 1)



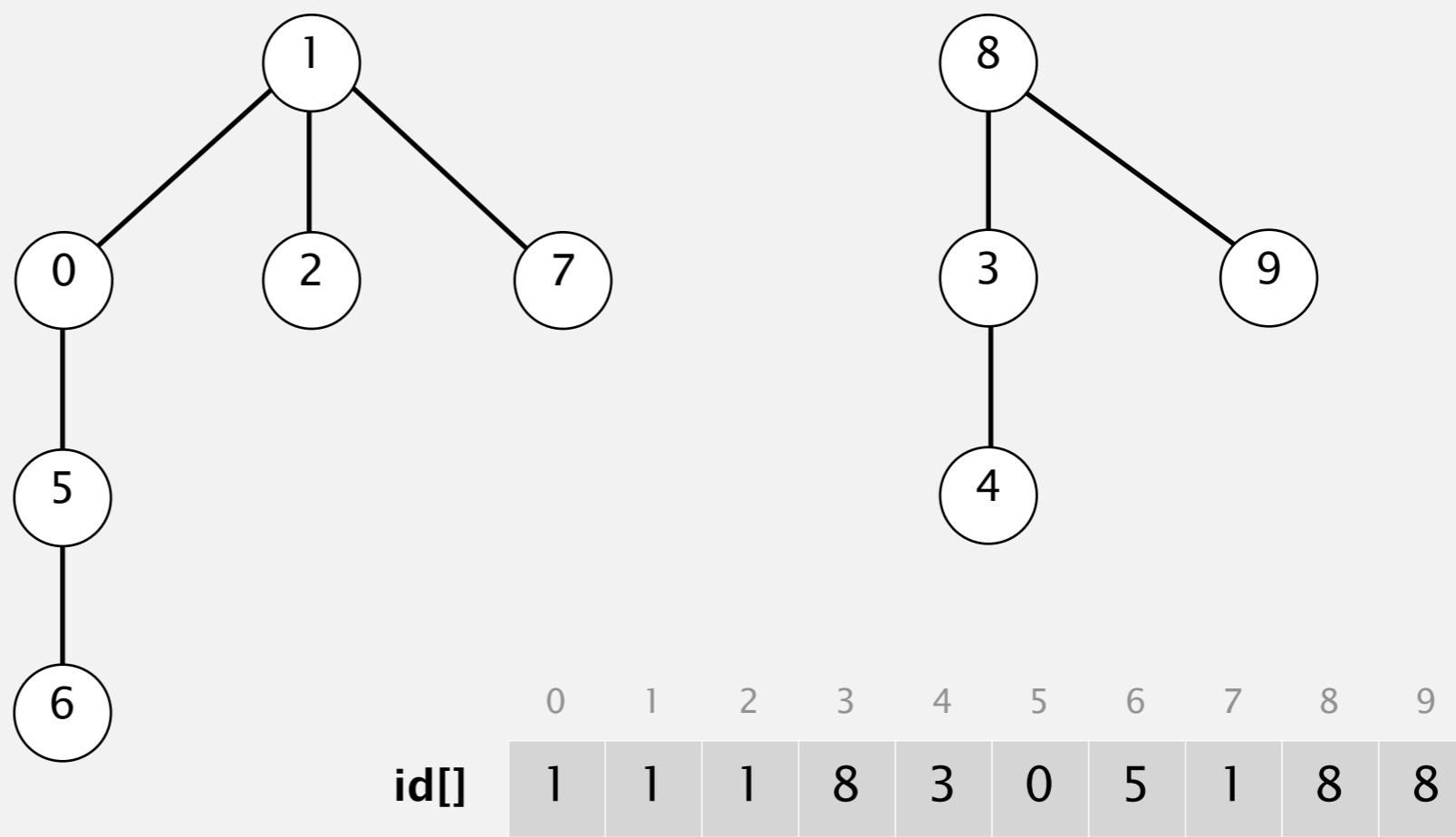
	0	1	2	3	4	5	6	7	8	9
id[]	0	1	1	8	3	0	5	1	8	8

Quick-union demo

union(6, 1)

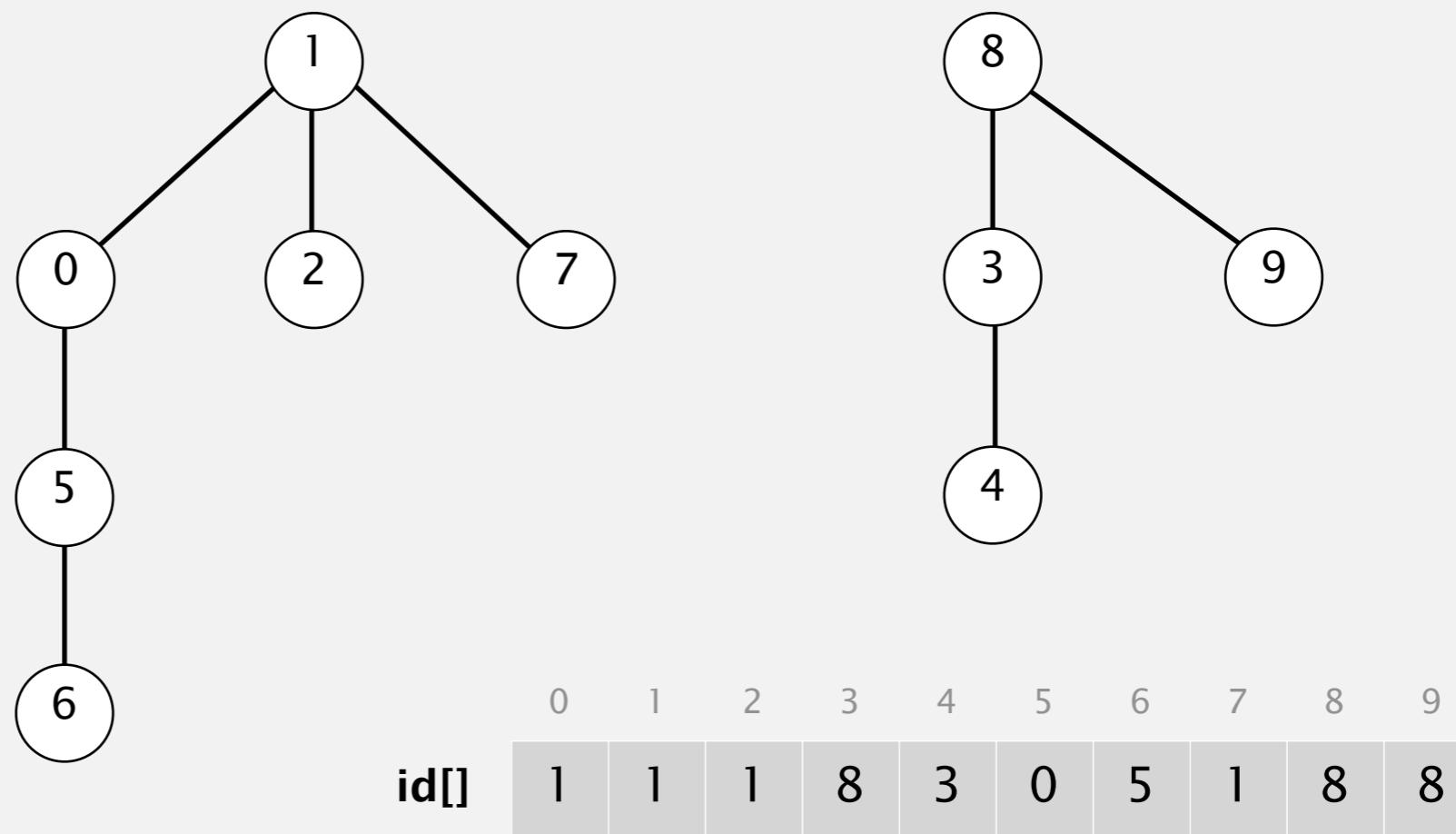


Quick-union demo



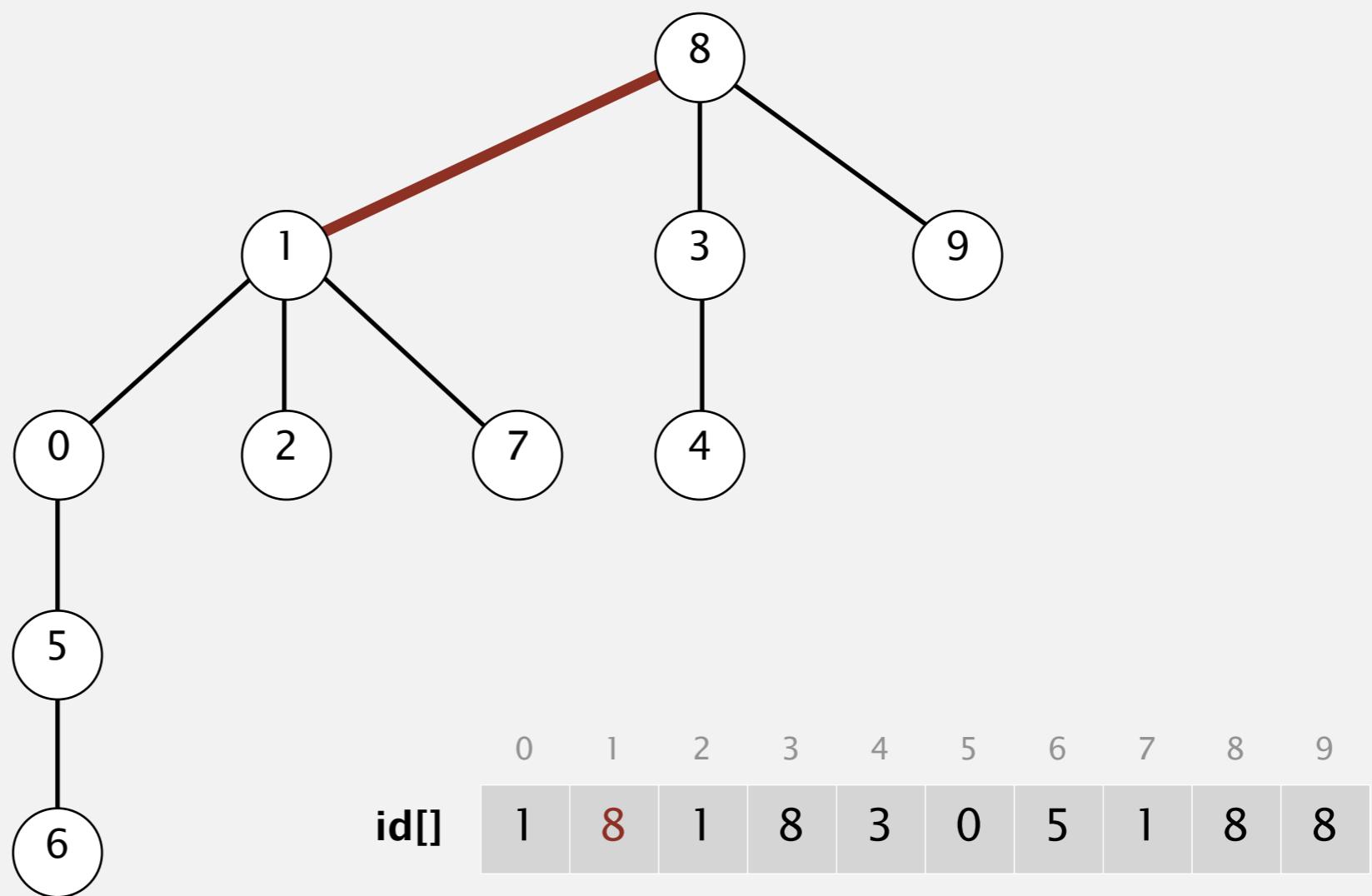
Quick-union demo

union(7, 3)

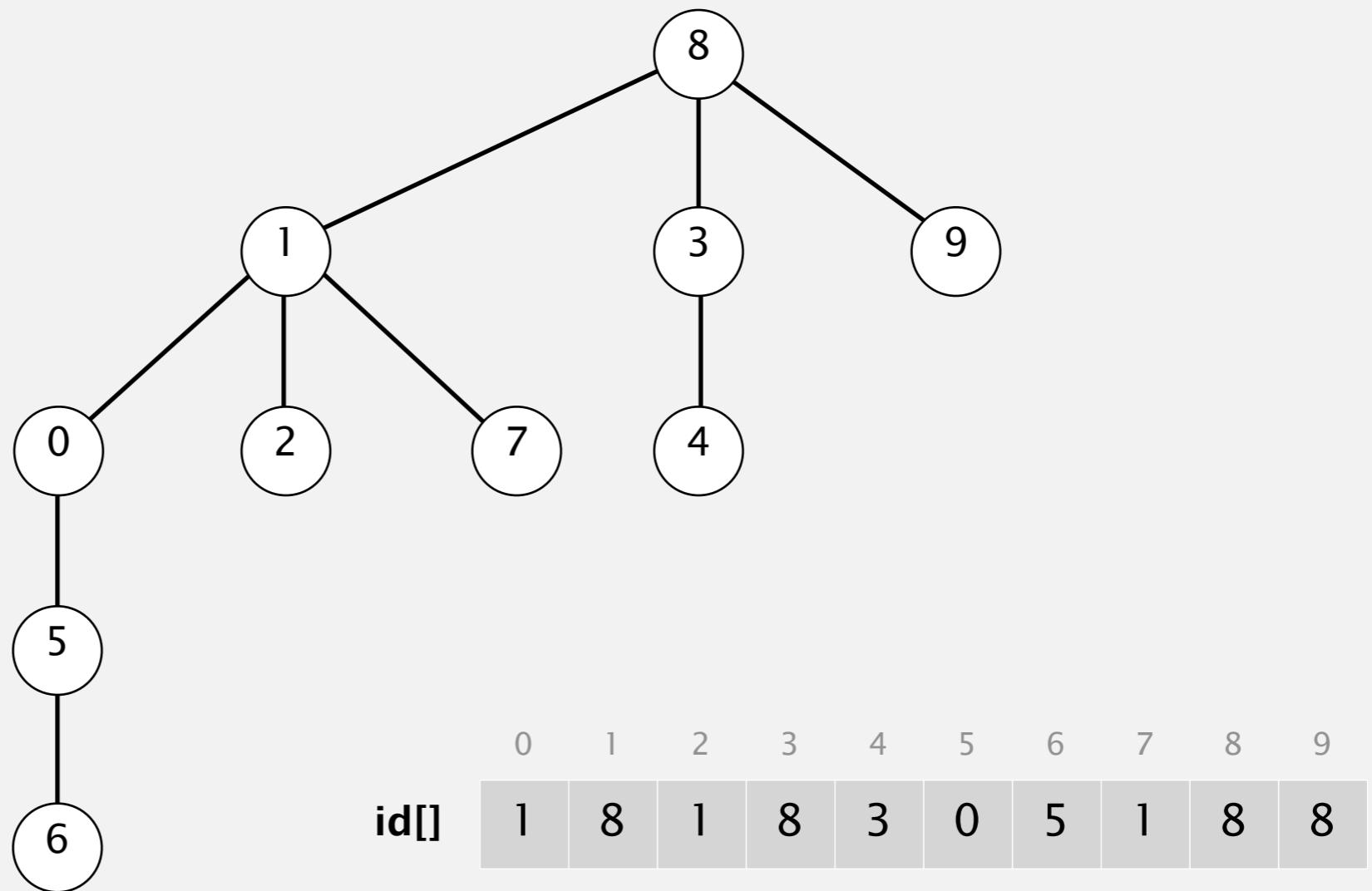


Quick-union demo

union(7, 3)



Quick-union demo



Quick-union: Java implementation

```
public class QuickUnionUF
{
    private int[] id;

    public QuickUnionUF(int N)
    {
        for (int i = 0; i < N; i++)
            id[i] = i; // set id of each object to itself
    }

    public int find(int i)
    {
        while (id[i] != i) // chase parent pointers until reach root
            i = id[i];
        return i;
    }

    public void union(int p, int q)
    {
        int pRoot = find(p);
        int qRoot = find(q);
        if (pRoot == qRoot)
            return;
        id[pRoot] = qRoot; // change root of p to point to root of q
    }
}
```

set id of each object to itself
(N array accesses)

chase parent pointers until reach root
(depth of i array accesses)

change root of p to point to root of q
(depth of p and q array accesses)

Quick-union: Java implementation

```
public class QuickUnionUF
{
    private int[] id;

    public QuickUnionUF(int N)
    {
        id = new int[N];
        for (int i = 0; i < N; i++) id[i] = i;
    }

    public int find(int i)
    {
        while (i != id[i]) i = id[i];
        return i;
    }

    public void union(int p, int q)
    {
        int i = find(p);
        int j = find(q);
        id[i] = j;
    }
}
```

set id of each object to itself
(N array accesses)

chase parent pointers until reach root
(depth of i array accesses)

change root of p to point to root of q
(depth of p and q array accesses)

Quick-union is also too slow

Cost model. Number of array accesses (for read or write).

algorithm	initialize	union	find	connected
quick-find	N	N	1	1
quick-union	N	N [†]	N	N

† includes cost of finding roots

← worst case

Quick-find defect.

- Union too expensive (N array accesses).
- Trees are flat, but too expensive to keep them flat.

Quick-union defect.

- Trees can get tall.
- Find/connected too expensive (could be N array accesses).

Algorithms

ROBERT SEDGEWICK | KEVIN WAYNE

<http://algs4.cs.princeton.edu>

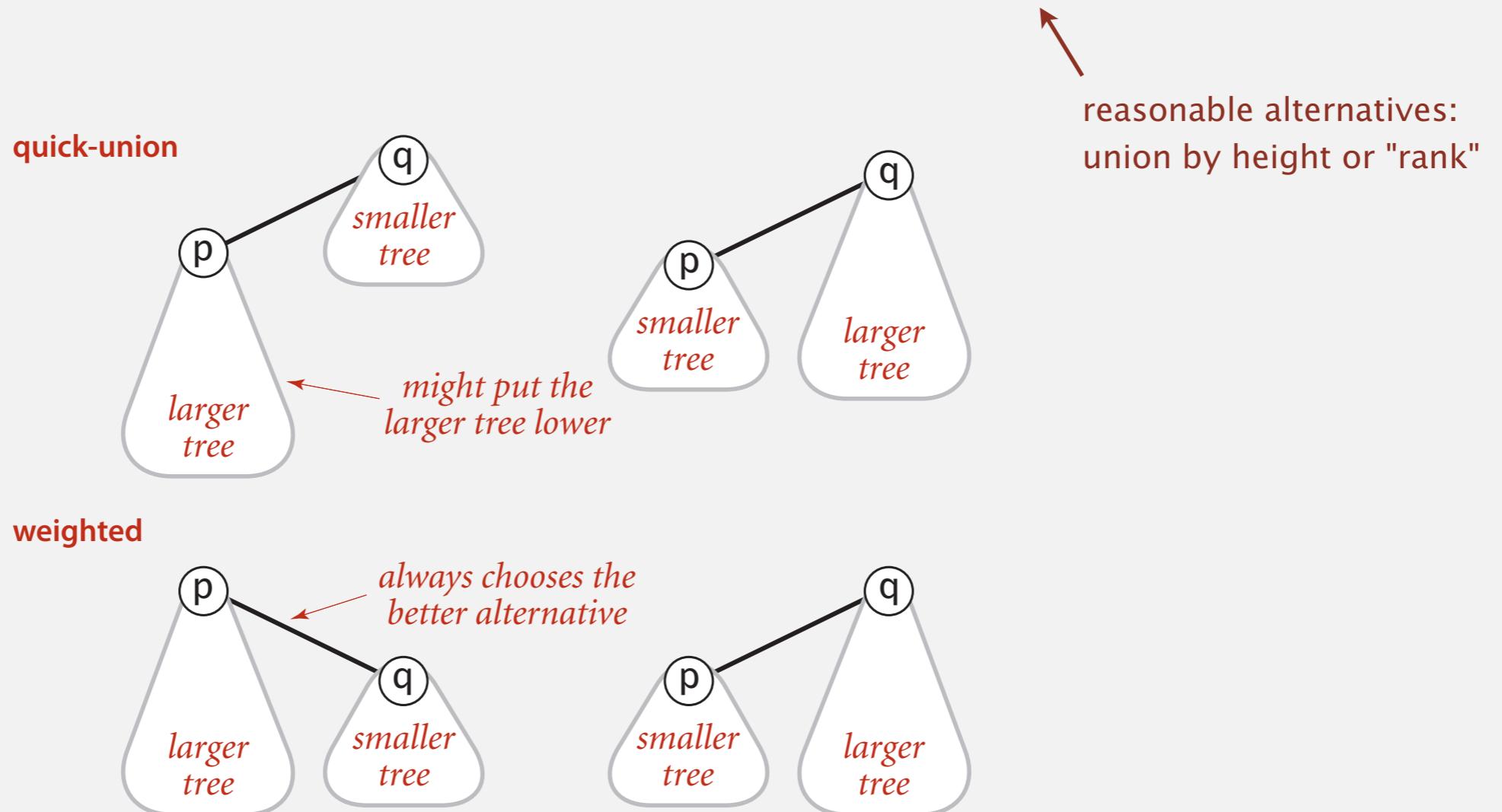
1.5 UNION-FIND

- ▶ *dynamic connectivity*
- ▶ *quick find*
- ▶ *quick union*
- ▶ ***improvements***
- ▶ *applications*

Improvement 1: weighting

Weighted quick-union.

- Modify quick-union to avoid tall trees.
- Keep track of size of each tree (number of objects).
- Balance by linking root of smaller tree to root of larger tree.



Weighted quick-union demo



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

Weighted quick-union demo



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

Weighted quick-union demo

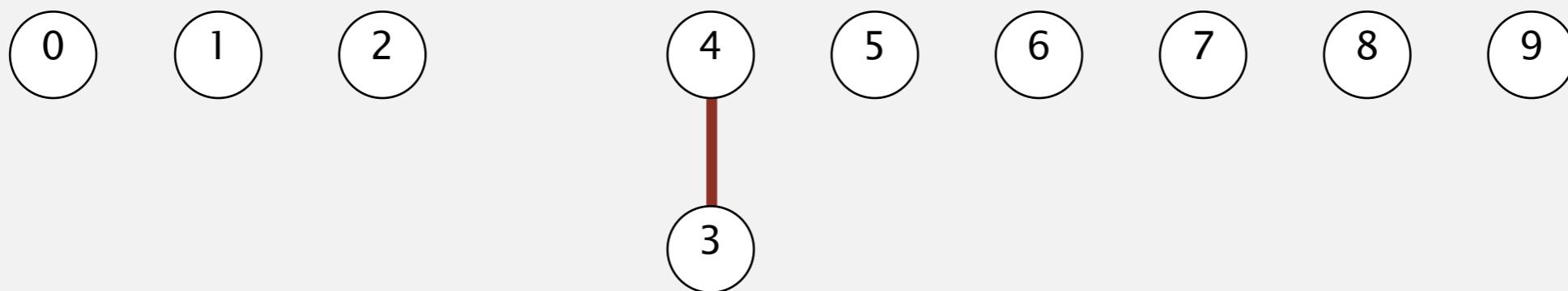
union(4, 3)



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

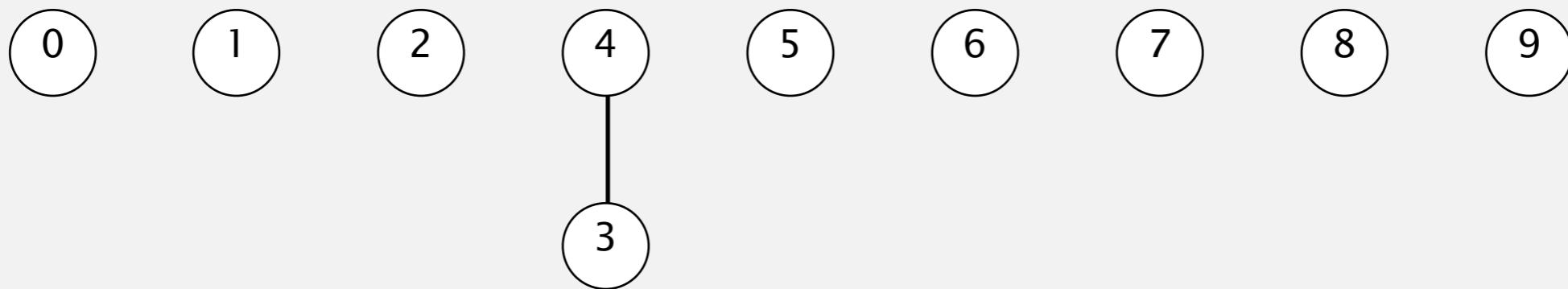
Weighted quick-union demo

union(4, 3)



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

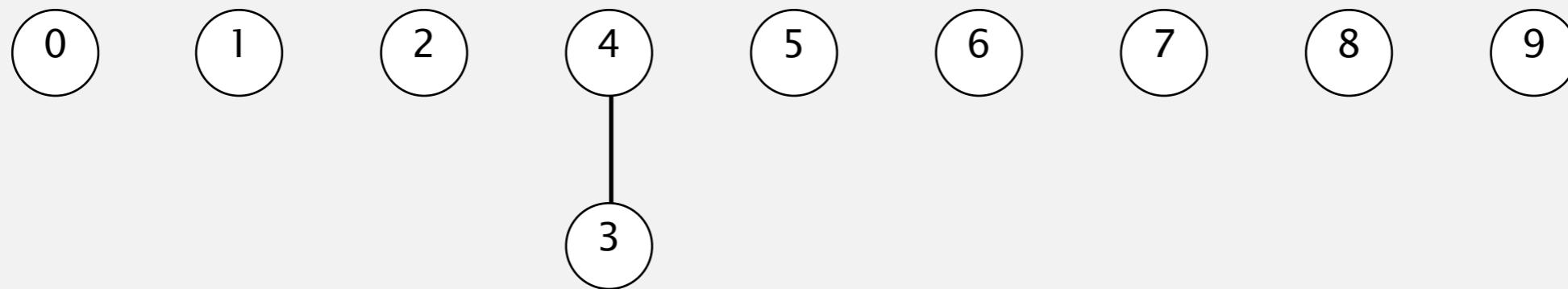
Weighted quick-union demo



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

Weighted quick-union demo

union(3, 8)

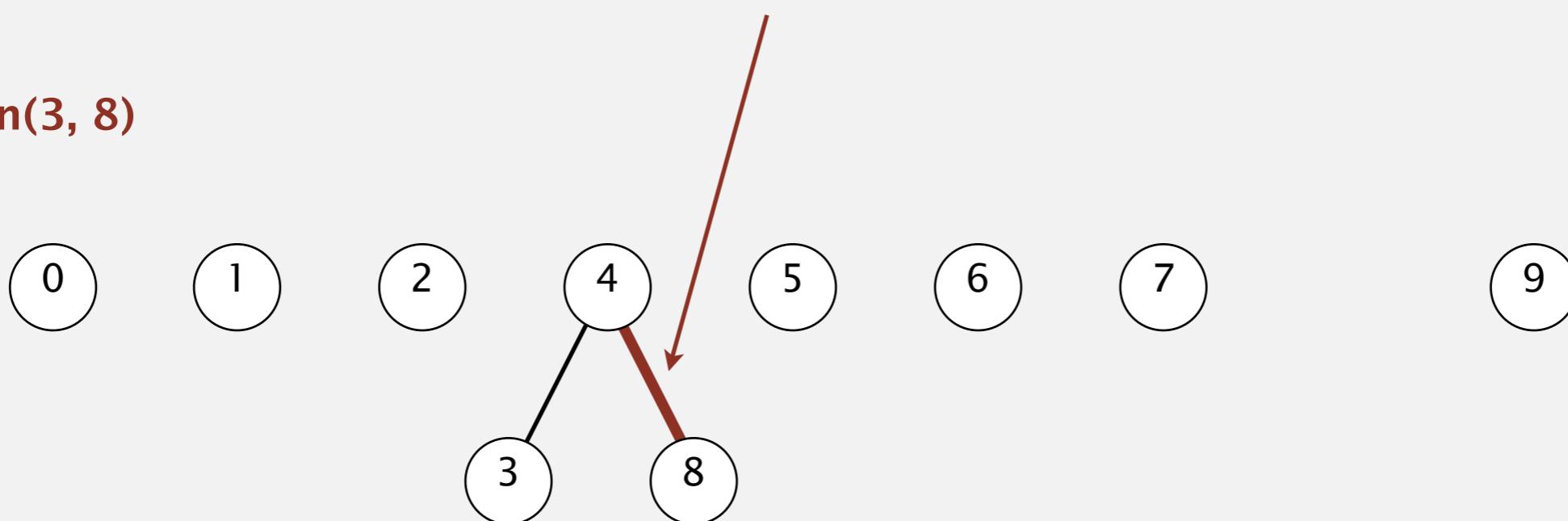


	0	1	2	3	4	5	6	7	8	9
id[]	0	1	2	4	4	5	6	7	8	9

Weighted quick-union demo

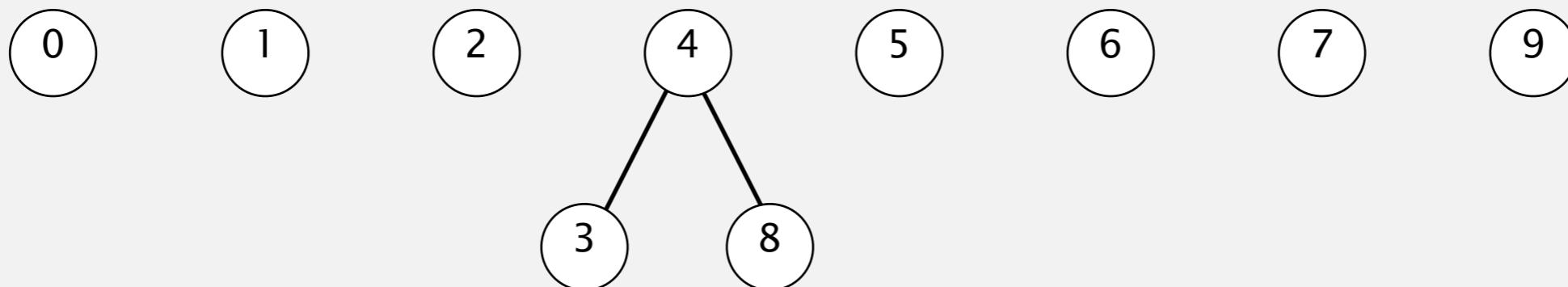
union(3, 8)

weighting: make 8 point to 4 (instead of 4 to 8)



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

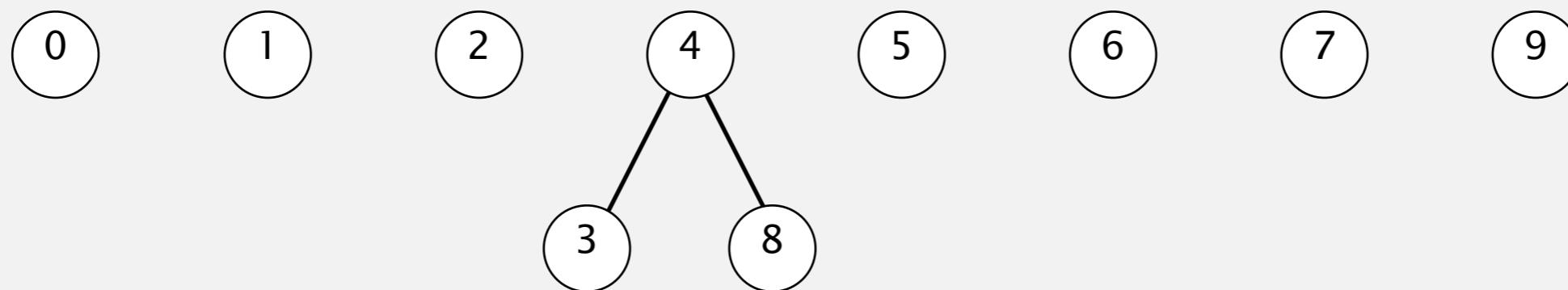
Weighted quick-union demo



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

Weighted quick-union demo

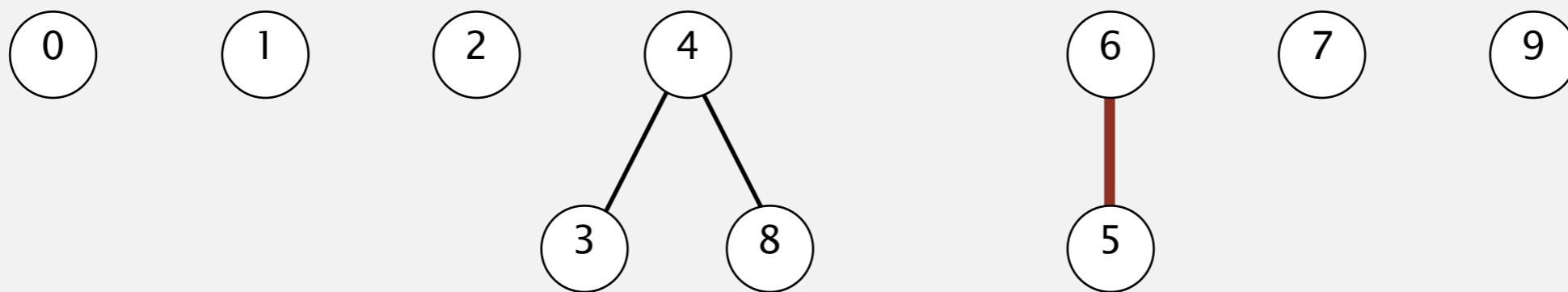
union(6, 5)



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

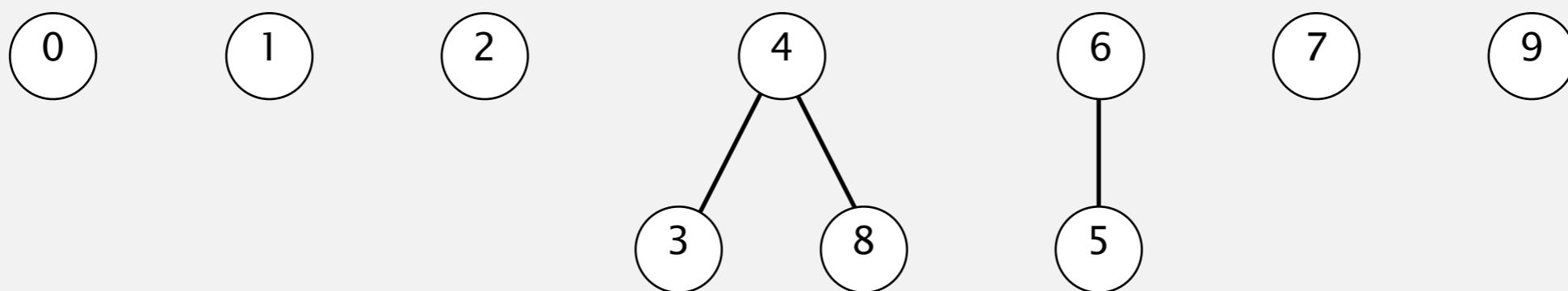
Weighted quick-union demo

union(6, 5)



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

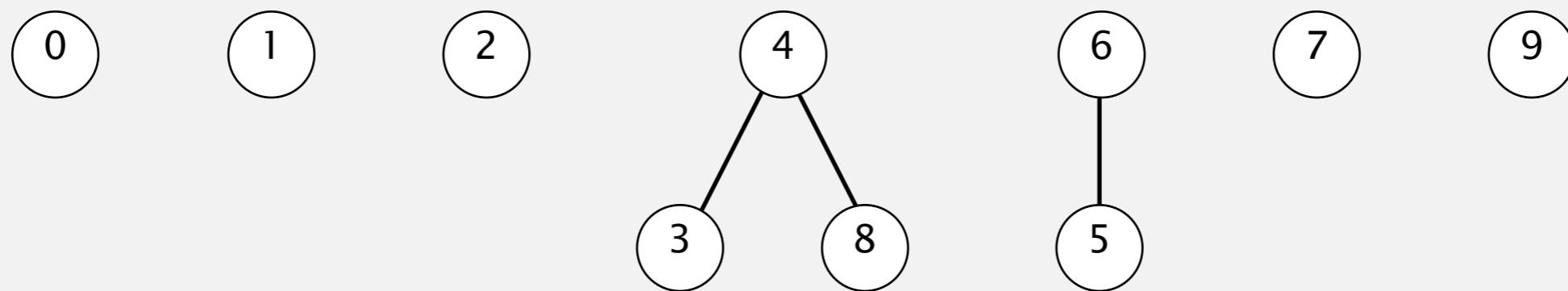
Weighted quick-union demo



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

Weighted quick-union demo

union(9, 4)

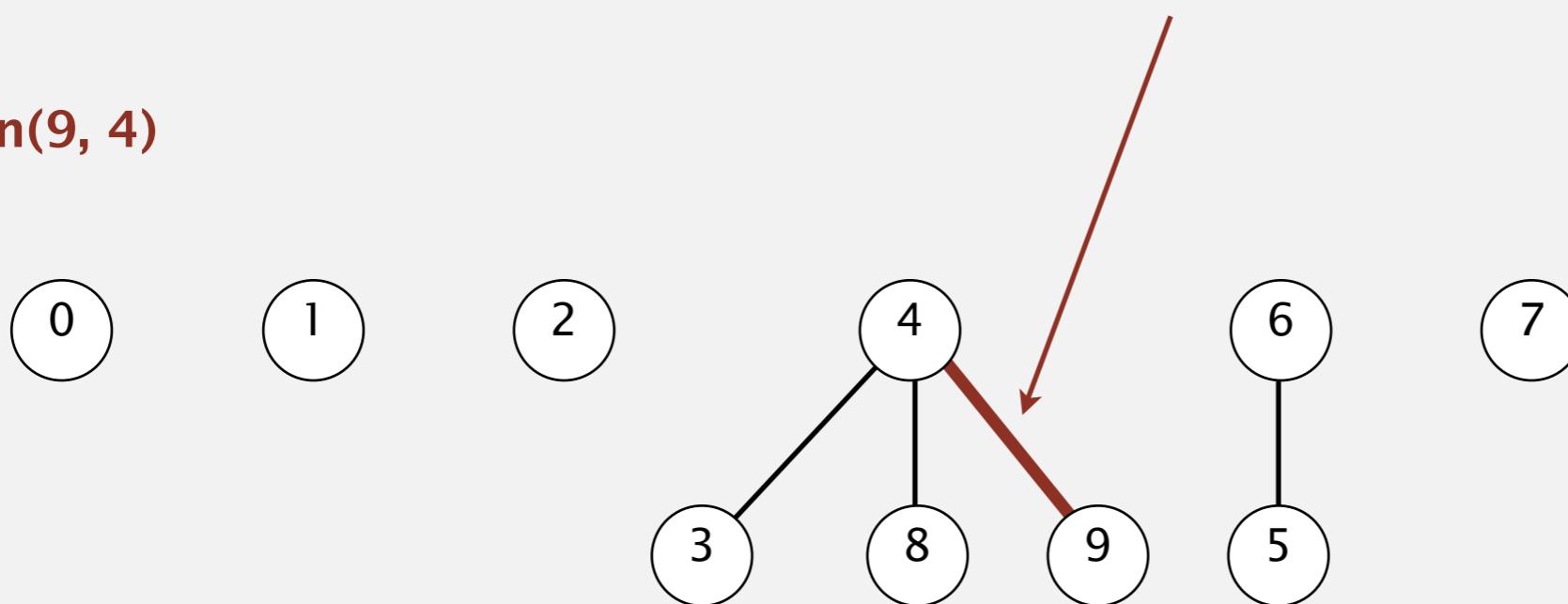


	0	1	2	3	4	5	6	7	8	9
id[]	0	1	2	4	4	6	6	7	4	9

Weighted quick-union demo

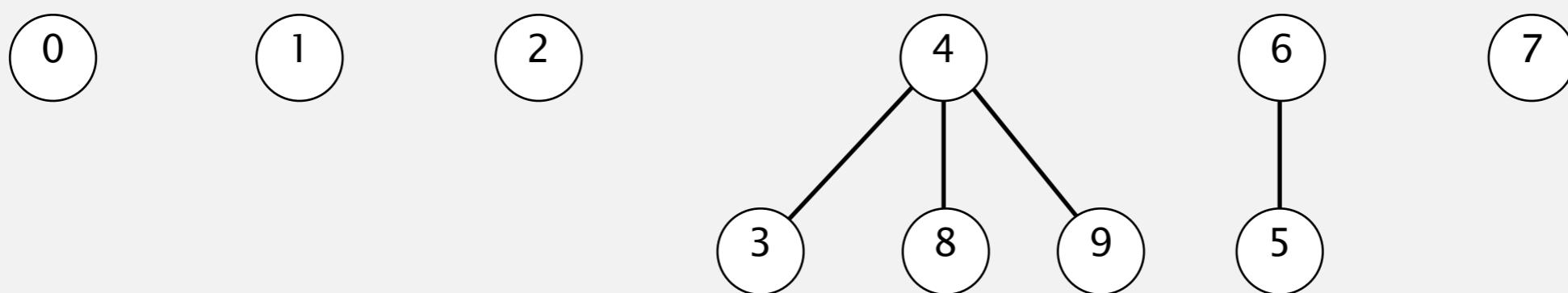
union(9, 4)

weighting: make 9 point to 4



0	1	2	3	4	4	6	6	7	4	4
id[]	0	1	2	4	4	6	6	7	4	4

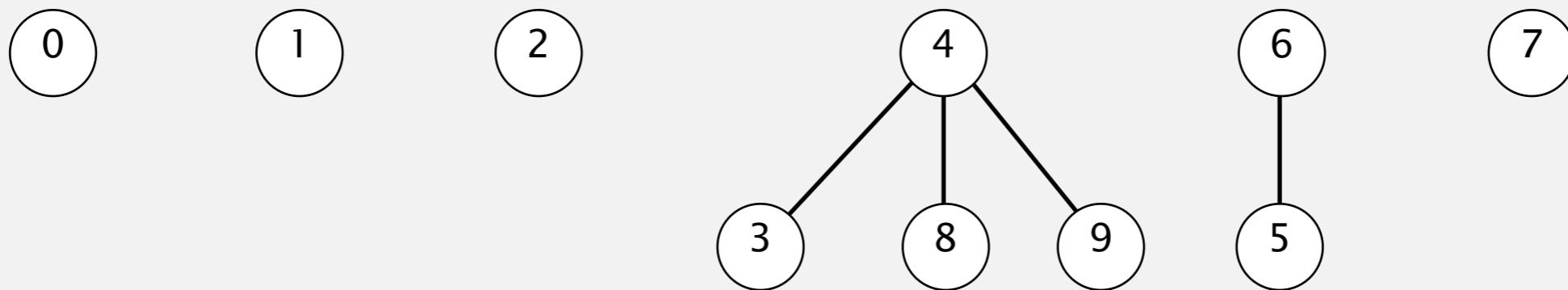
Weighted quick-union demo



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

Weighted quick-union demo

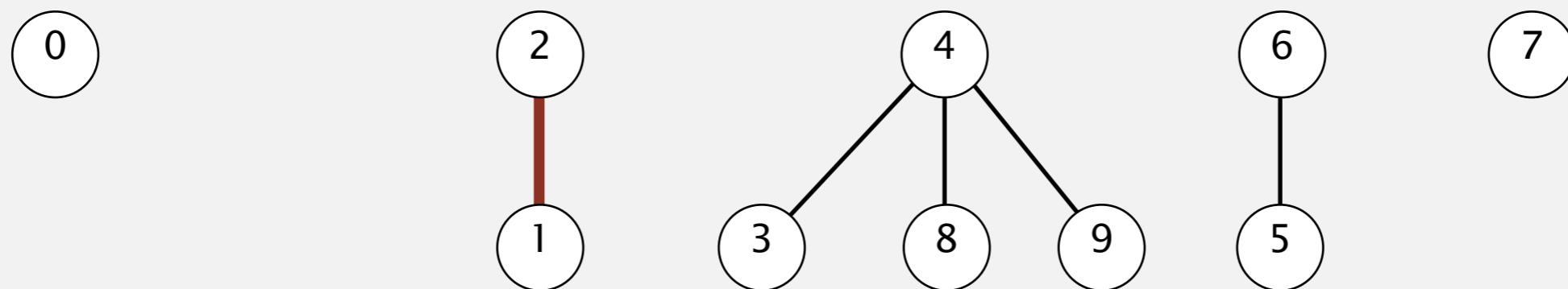
union(2, 1)



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

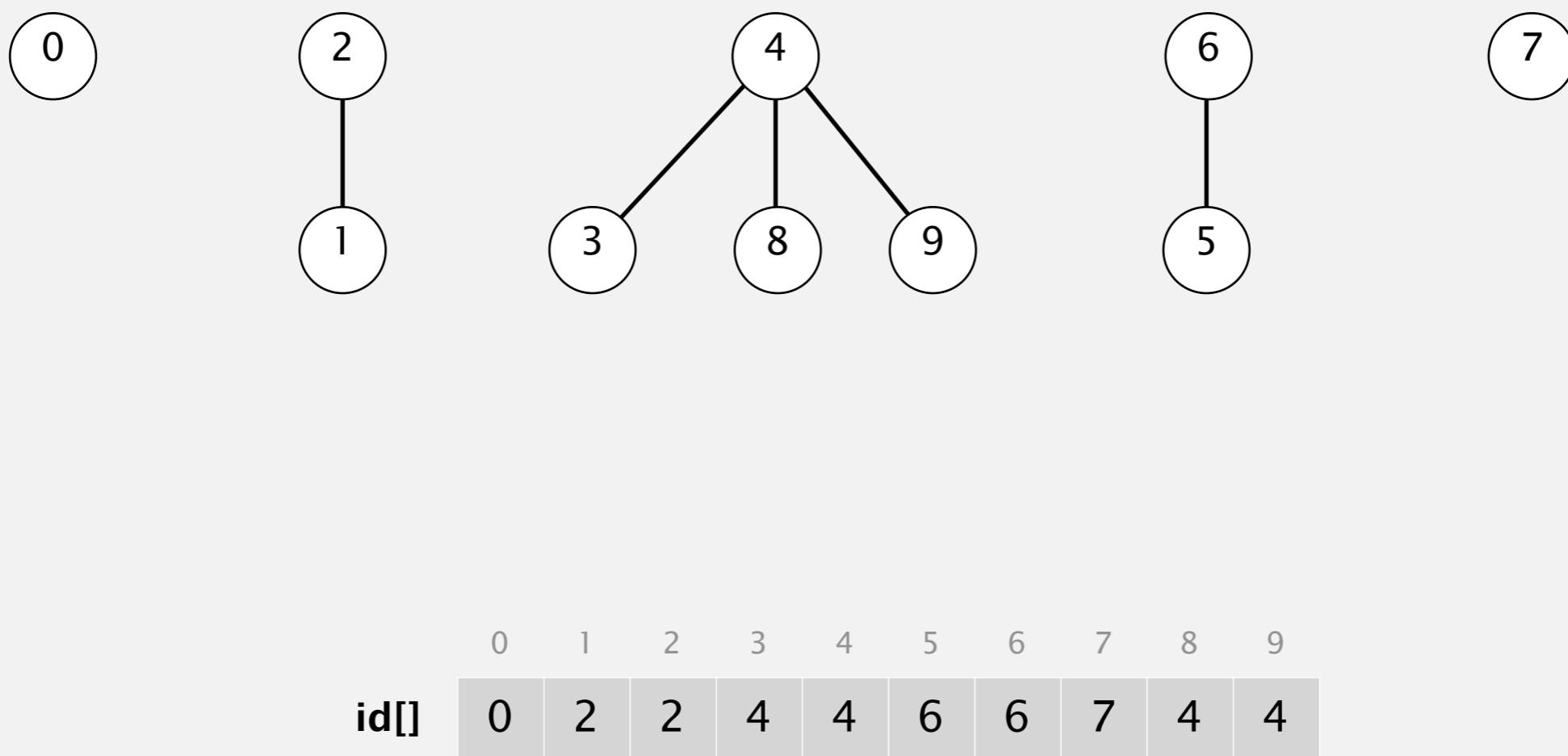
Weighted quick-union demo

union(2, 1)



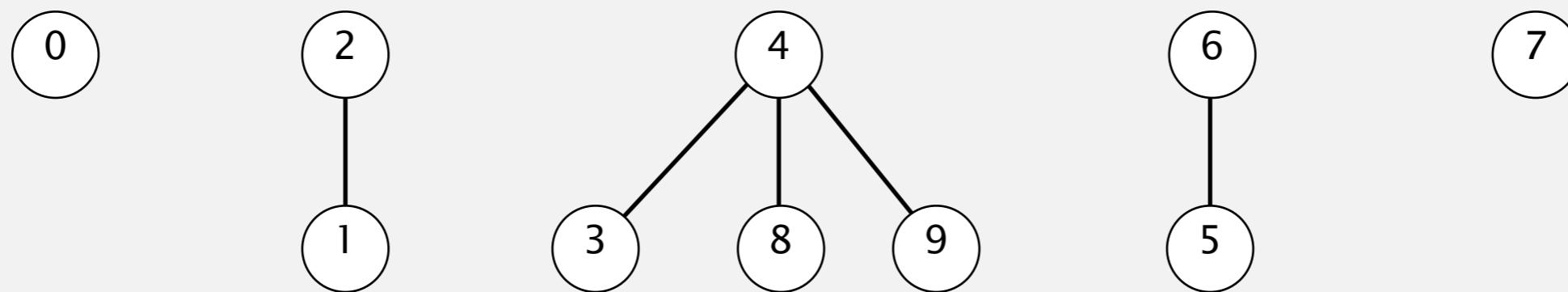
0	1	2	3	4	5	6	7	8	9	
id[]	0	2	2	4	4	6	6	7	4	4

Weighted quick-union demo



Weighted quick-union demo

union(5, 0)

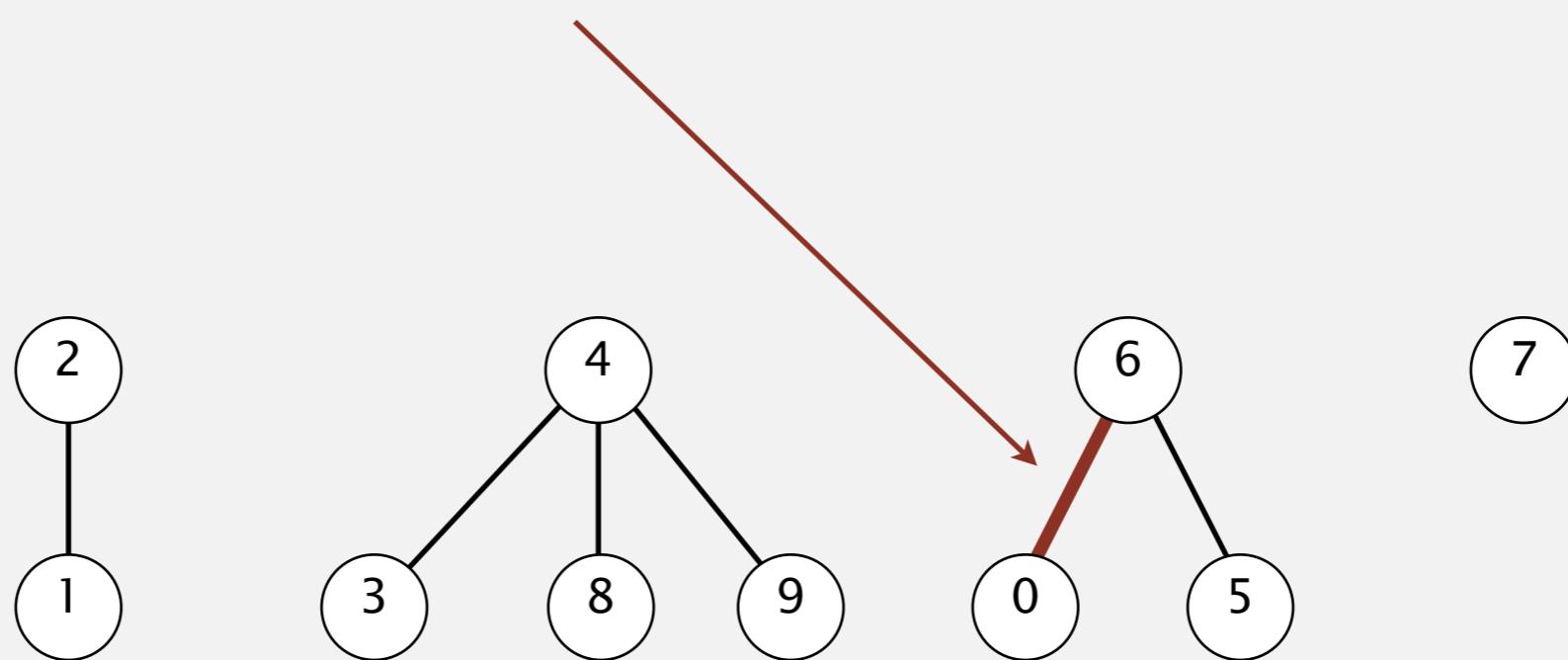


	0	1	2	3	4	5	6	7	8	9
id[]	0	2	2	4	4	6	6	7	4	4

Weighted quick-union demo

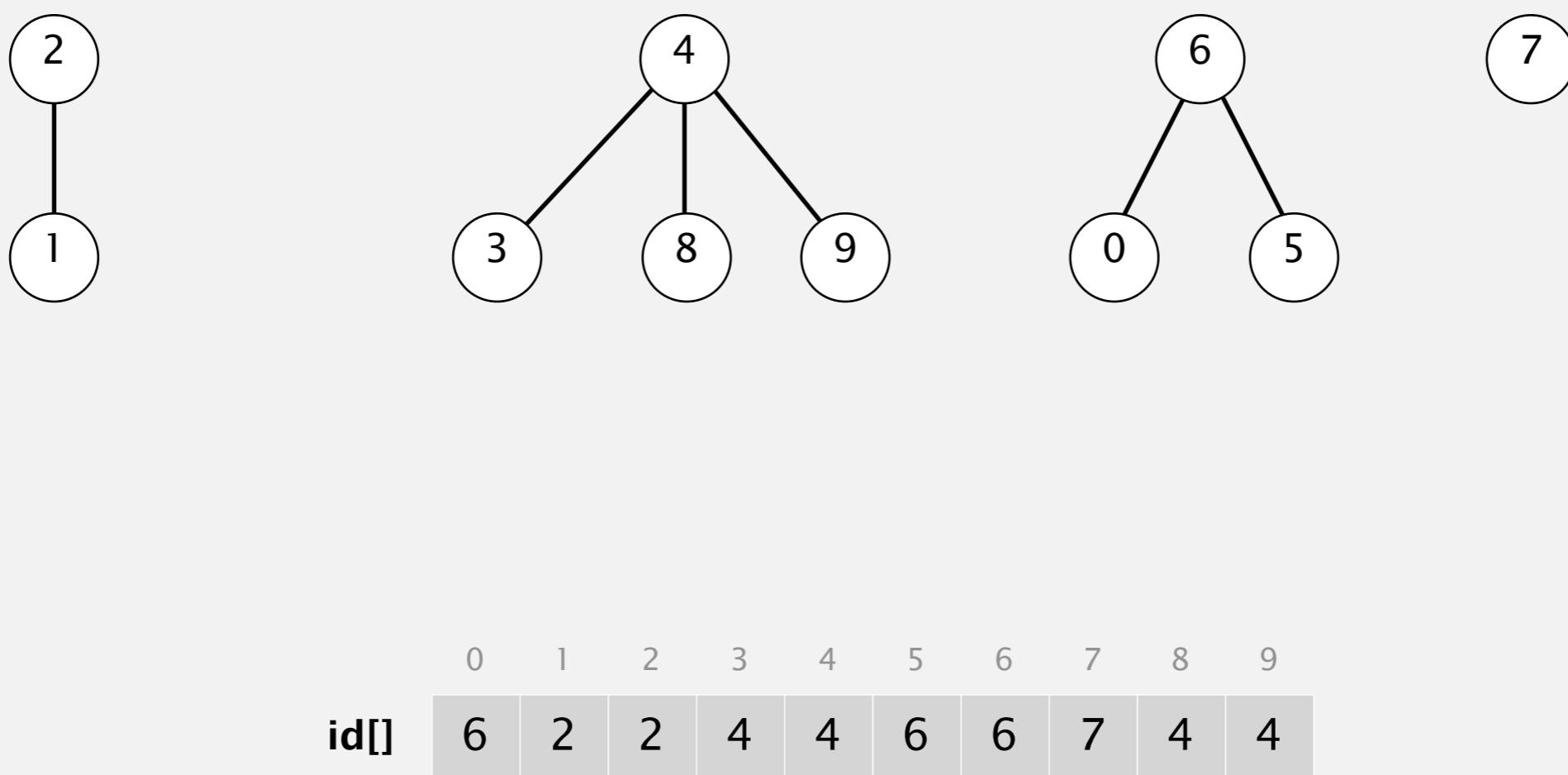
union(5, 0)

weighting: make 0 point to 6 (instead of 6 to 0)



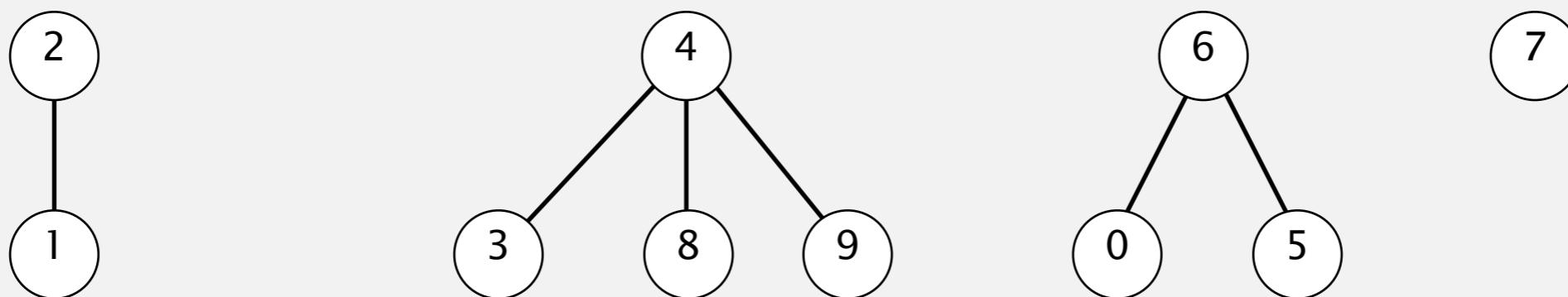
0	1	2	3	4	5	6	7	8	9	
id[]	6	2	2	4	4	6	6	7	4	4

Weighted quick-union demo



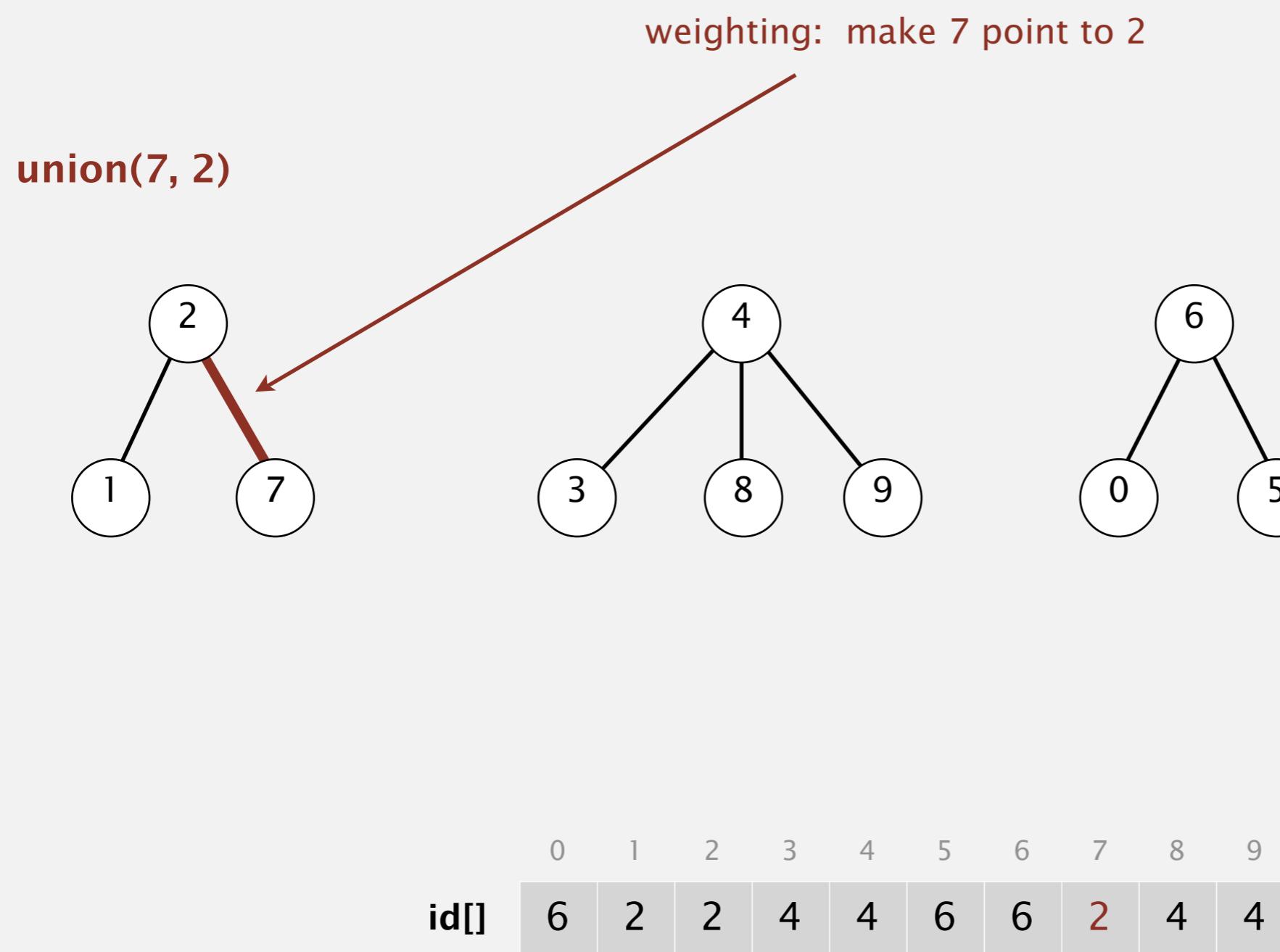
Weighted quick-union demo

union(7, 2)

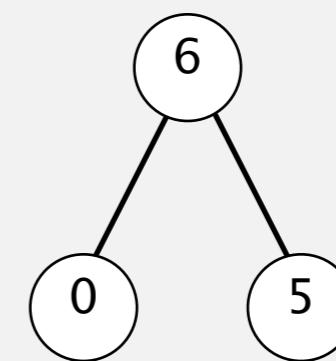
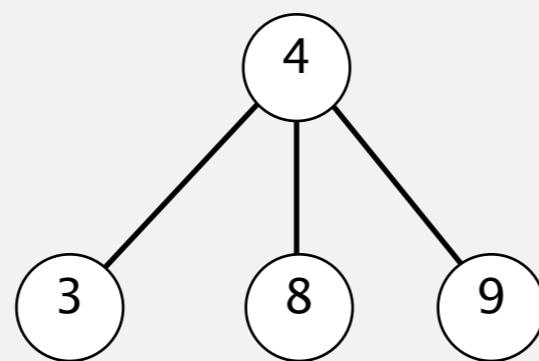
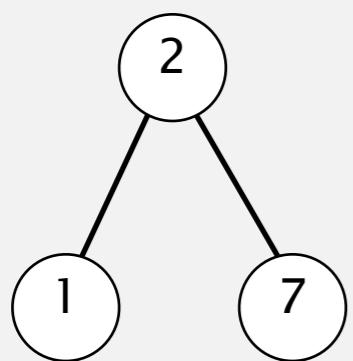


0	1	2	3	4	5	6	7	8	9	
id[]	6	2	2	4	4	6	6	7	4	4

Weighted quick-union demo



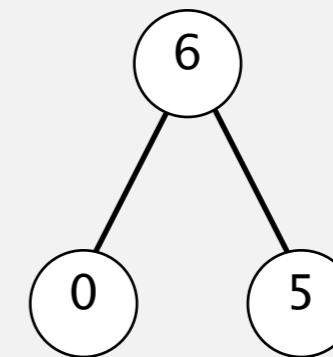
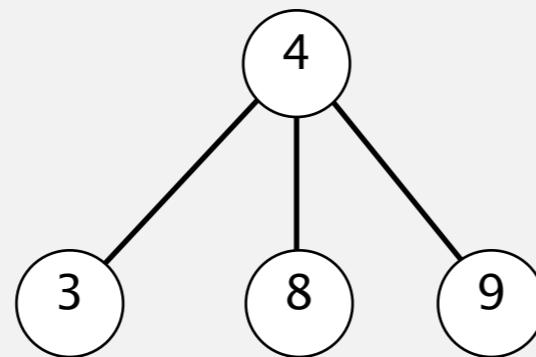
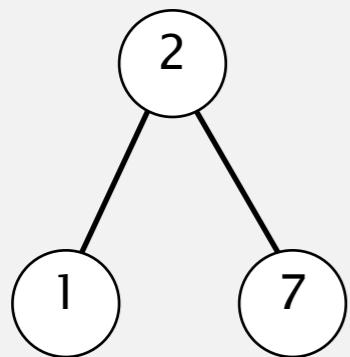
Weighted quick-union demo



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

Weighted quick-union demo

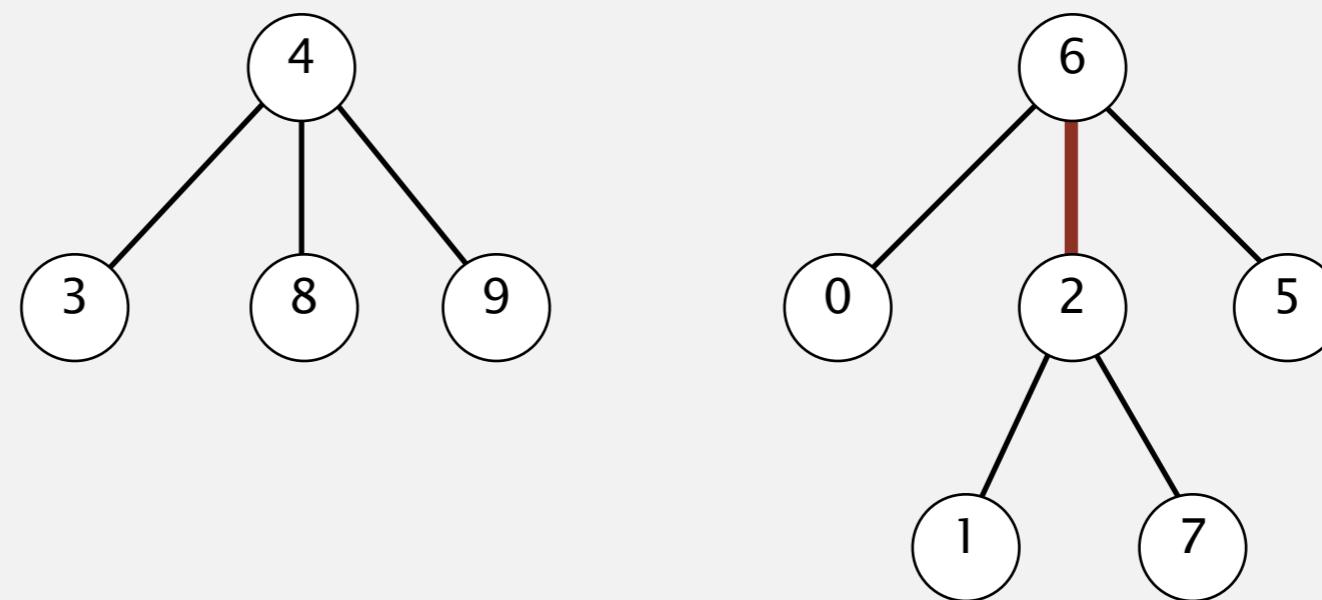
union(6, 1)



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

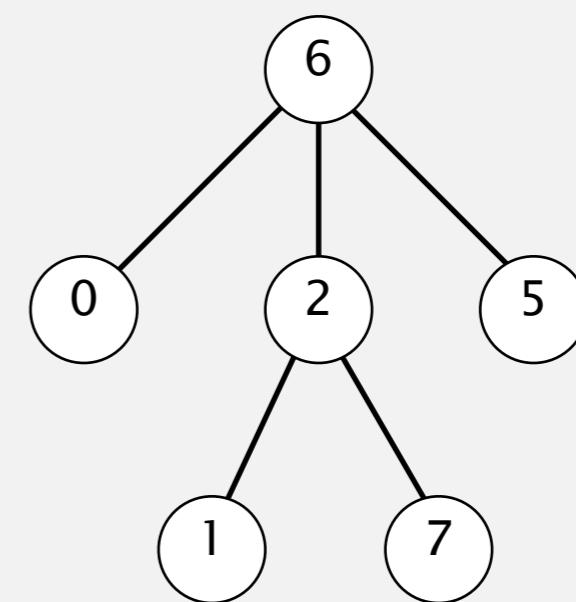
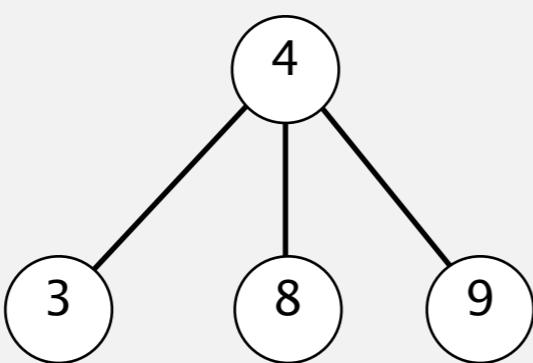
Weighted quick-union demo

union(6, 1)



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

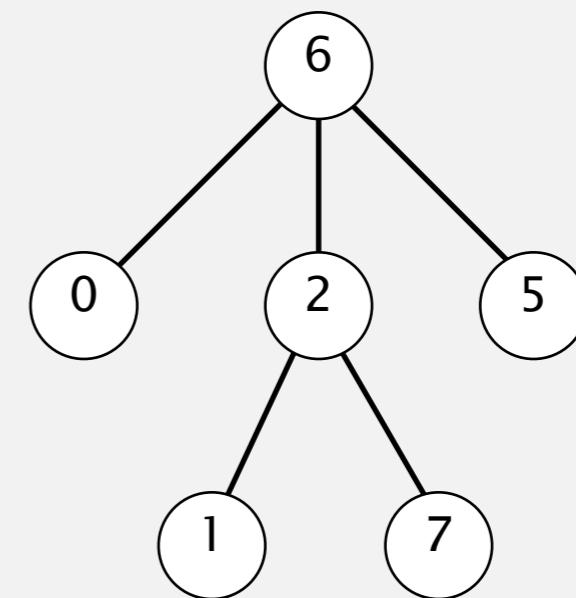
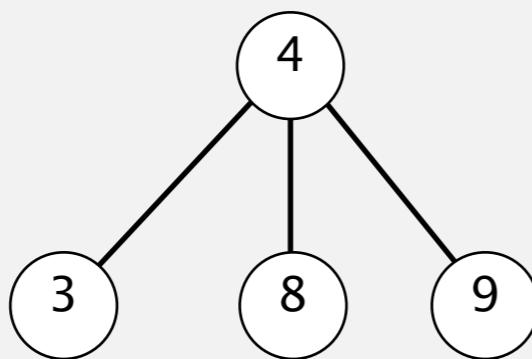
Weighted quick-union demo



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

Weighted quick-union demo

union(7, 3)

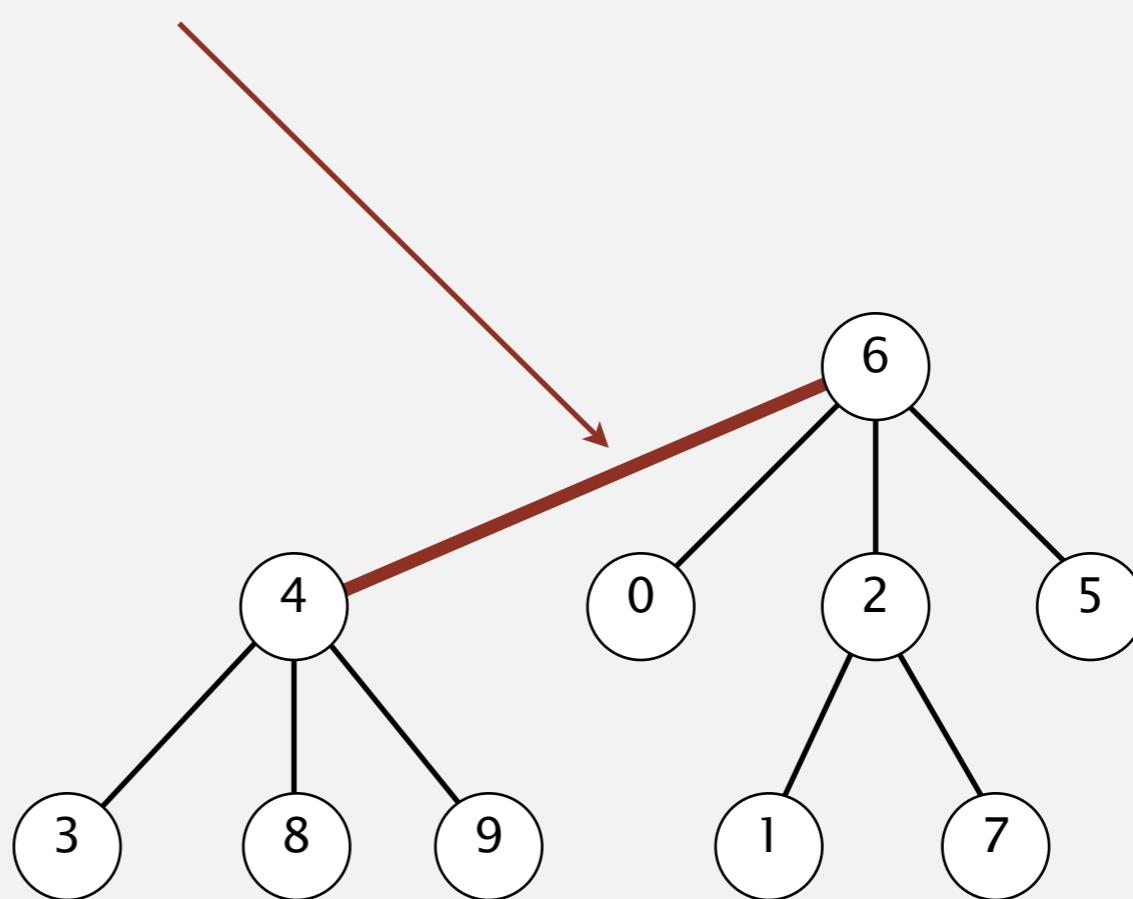


0	1	2	3	4	5	6	7	8	9	
id[]	6	2	6	4	4	6	6	2	4	4

Weighted quick-union demo

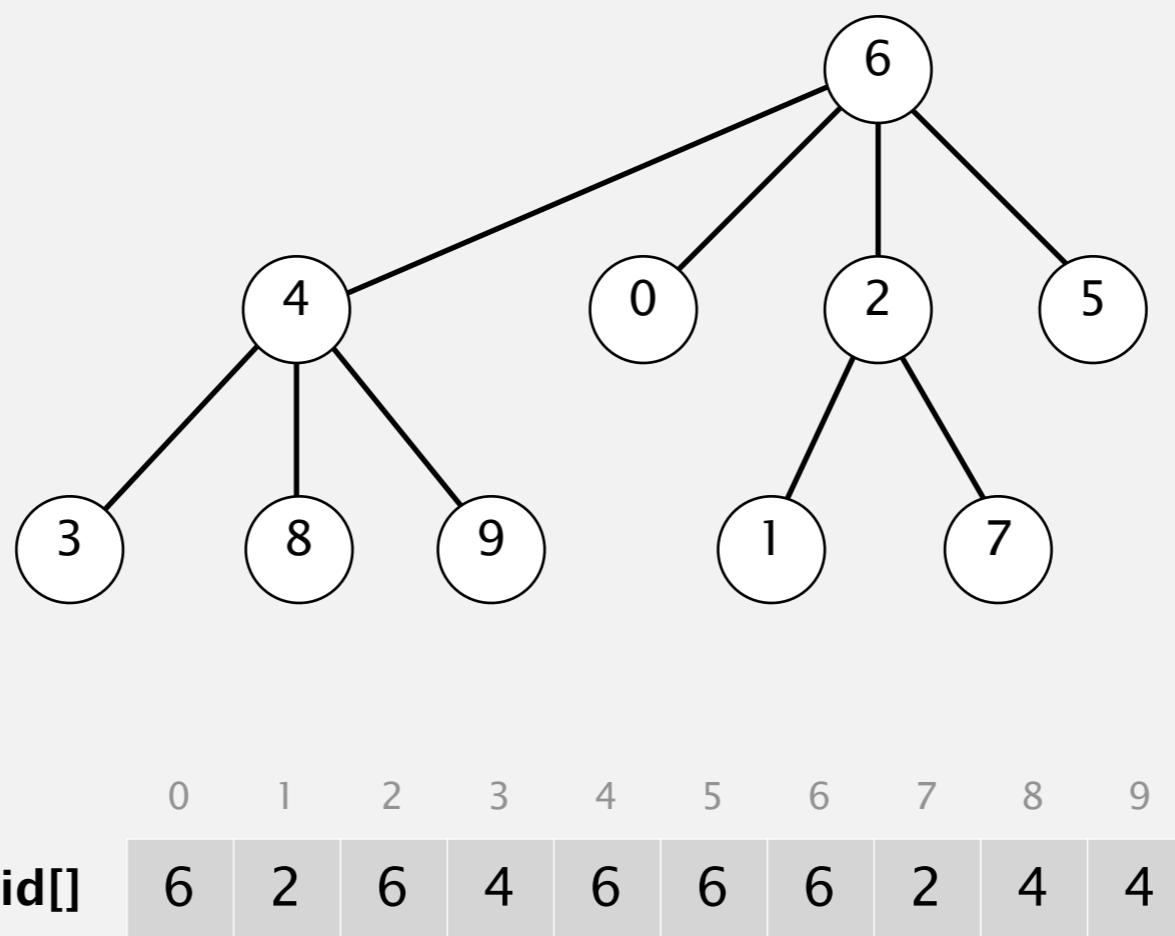
union(7, 3)

weighting: make 4 point to 6 (instead of 6 to 4)



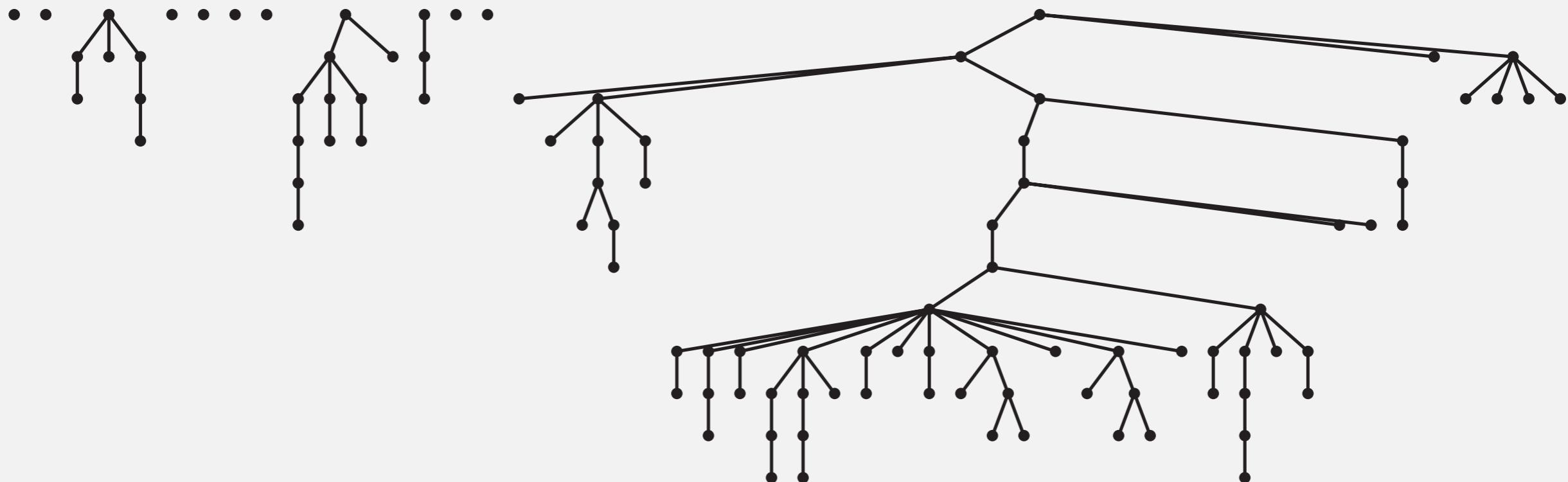
0	1	2	3	4	5	6	7	8	9	
id[]	6	2	6	4	6	6	6	2	4	4

Weighted quick-union demo

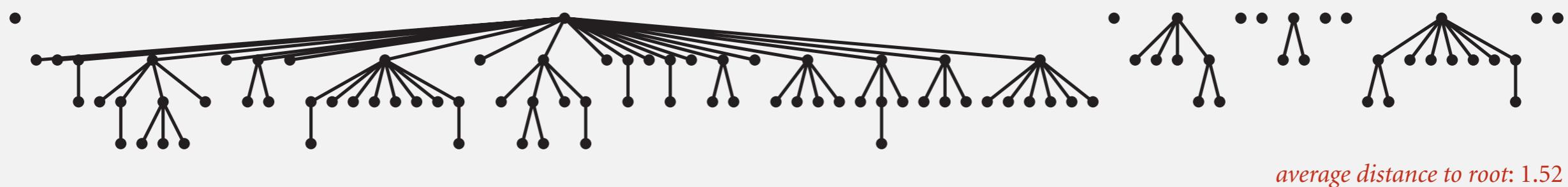


Example trees in quick-union versus weighted quick-union

quick-union



weighted



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

Weighted quick-union: Java implementation

Data structure. Same as quick-union, but maintain extra array $sz[i]$ to count number of objects in the tree rooted at i .

Find/connected. Identical to quick-union.

Union. Modify quick-union to:

- Link root of smaller tree to root of larger tree.
- Update the $sz[]$ array.

```
int i = find(p);
int j = find(q);
if (i == j) return;
if (sz[i] < sz[j]) { id[i] = j; sz[j] += sz[i]; }
else                  { id[j] = i; sz[i] += sz[j]; }
```

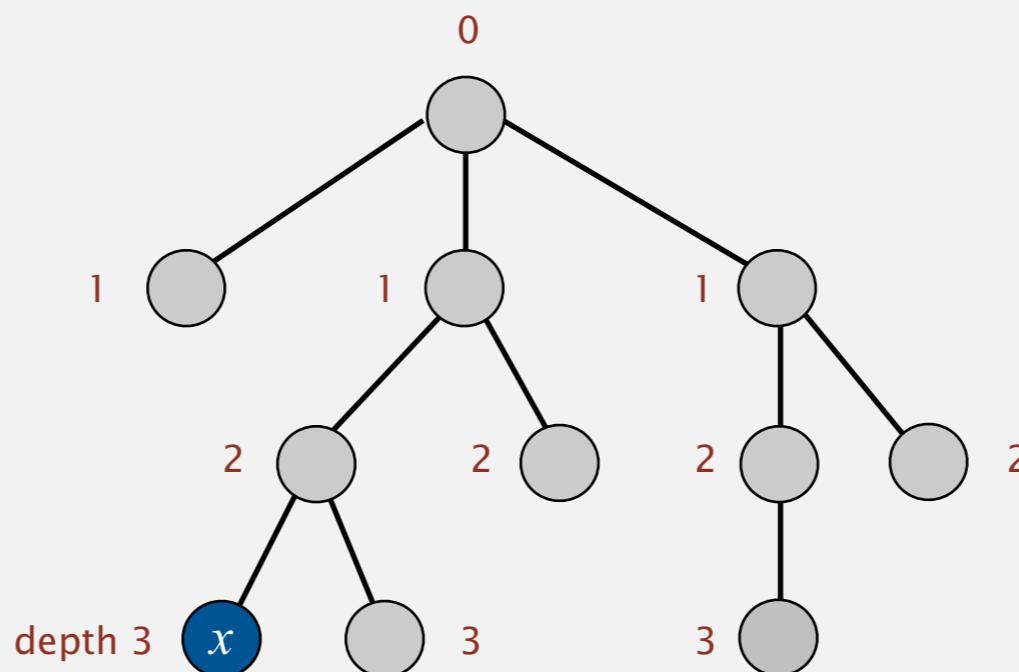
Weighted quick-union analysis

Running time.

- Find: takes time proportional to depth of p .
 - Union: takes constant time, given roots.

Proposition. Depth of any node x is at most $\lg N$.

\lg = base-2 logarithm



$$N = 11$$

$$\text{depth}(x) = 3 \leq \lg N$$

Weighted quick-union analysis

Running time.

- Find: takes time proportional to depth of p .
- Union: takes constant time, given roots.

\lg = base-2 logarithm

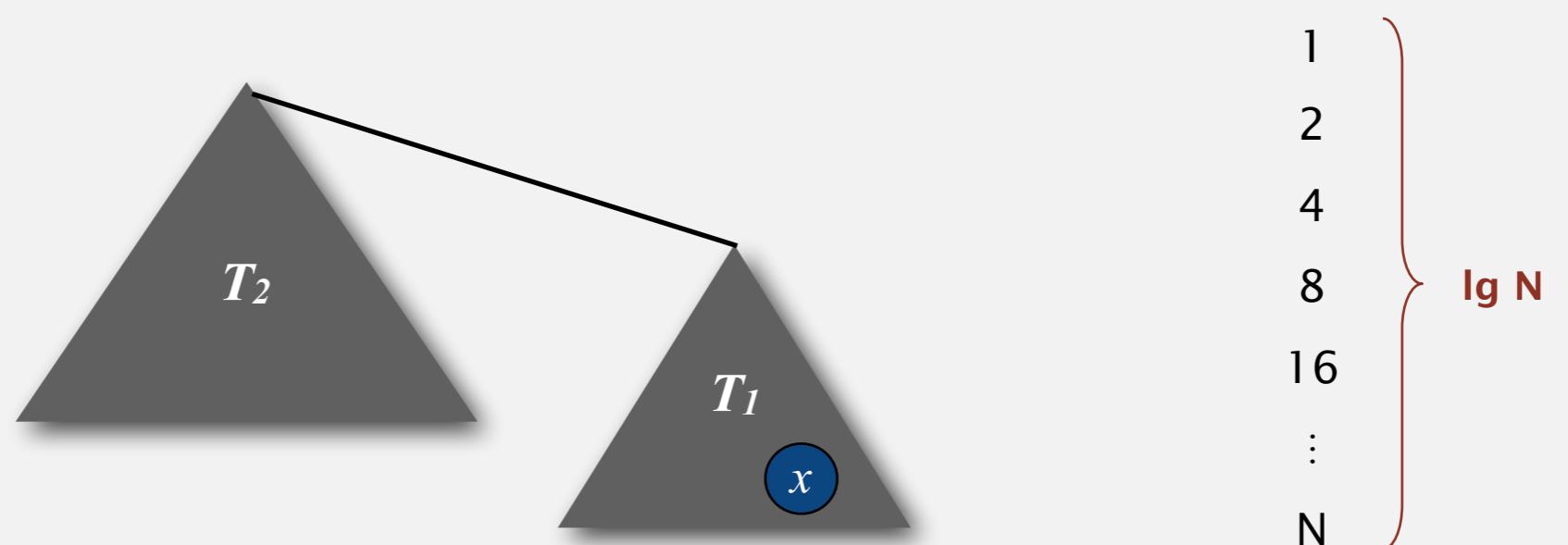


Proposition. Depth of any node x is at most $\lg N$.

Pf. What causes the depth of object x to increase?

Increases by 1 when tree T_1 containing x is merged into another tree T_2 .

- The size of the new tree containing x at least doubles since $|T_2| \geq |T_1|$.
- Size of tree containing x can double at most $\lg N$ times. Why?



Weighted quick-union analysis

Running time.

- Find: takes time proportional to depth of p .
- Union: takes constant time, given roots.

Proposition. Depth of any node x is at most $\lg N$.

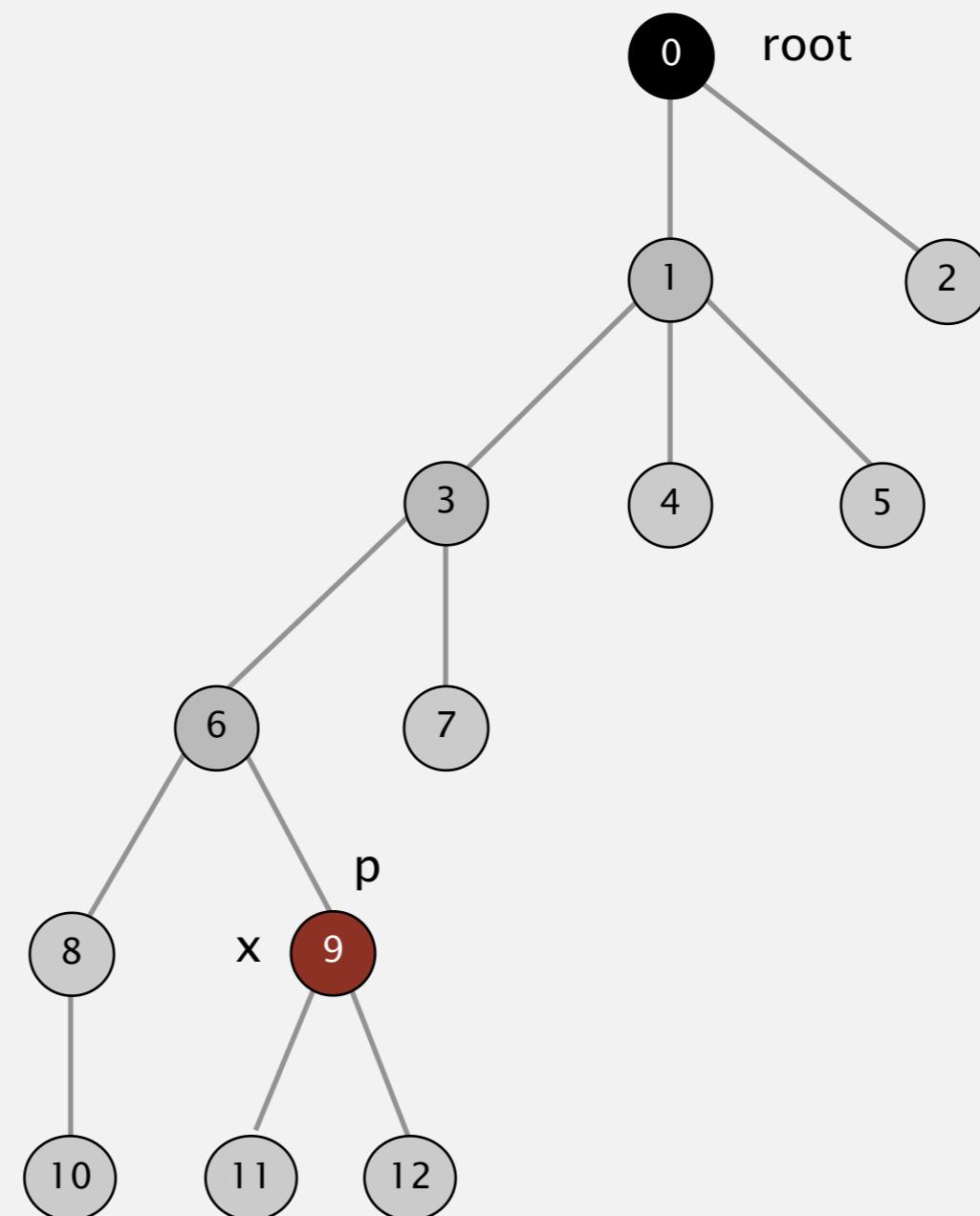
algorithm	initialize	union	find	connected
quick-find	N	N	1	1
quick-union	N	N^{\dagger}	N	N
weighted QU	N	$\lg N^{\dagger}$	$\lg N$	$\lg N$

\dagger includes cost of finding roots

- Q.** Stop at guaranteed acceptable performance?
A. No, easy to improve further.

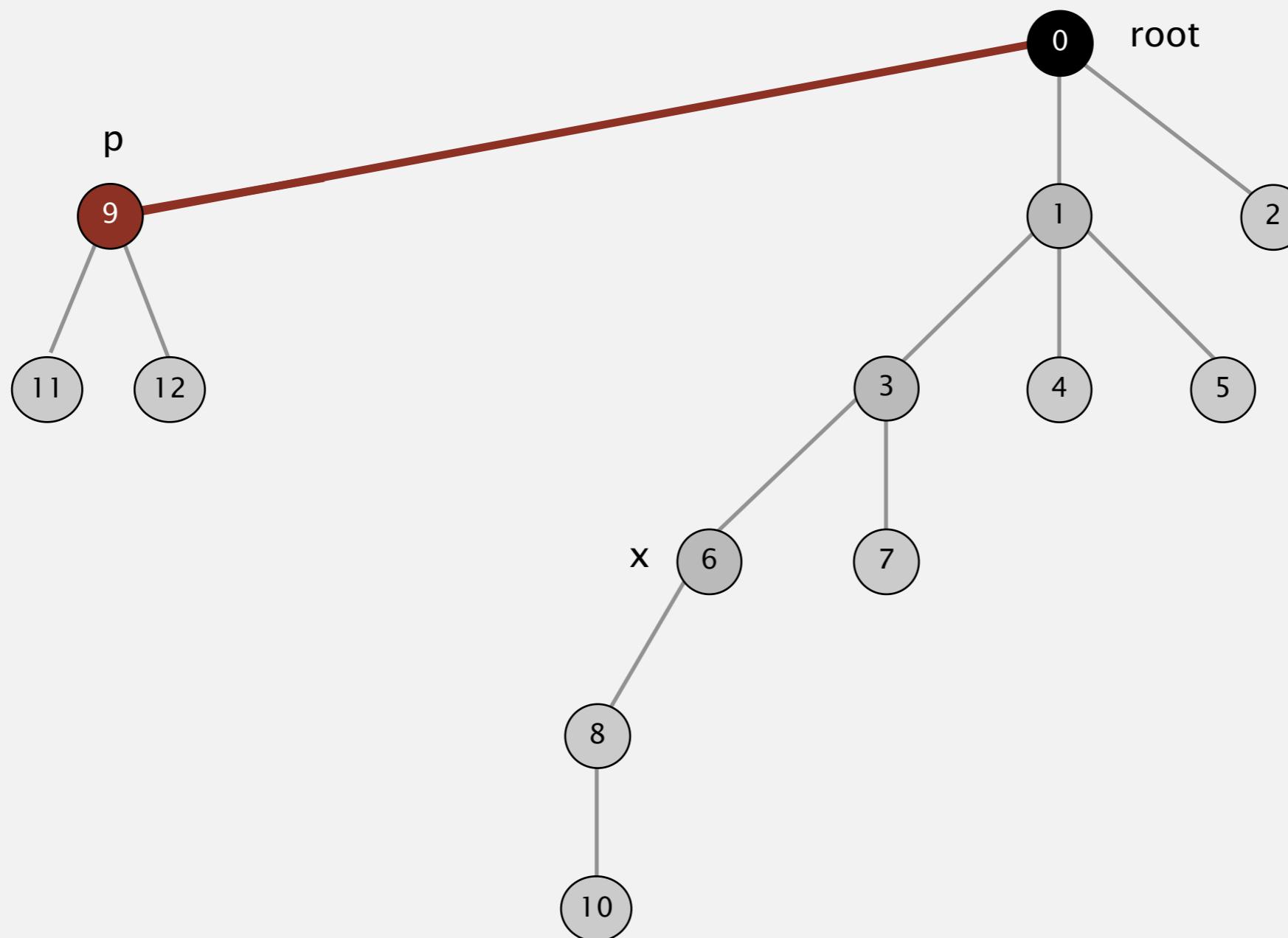
Improvement 2: path compression

Quick union with path compression. Just after computing the root of p , set the `id[]` of each examined node to point to that root.



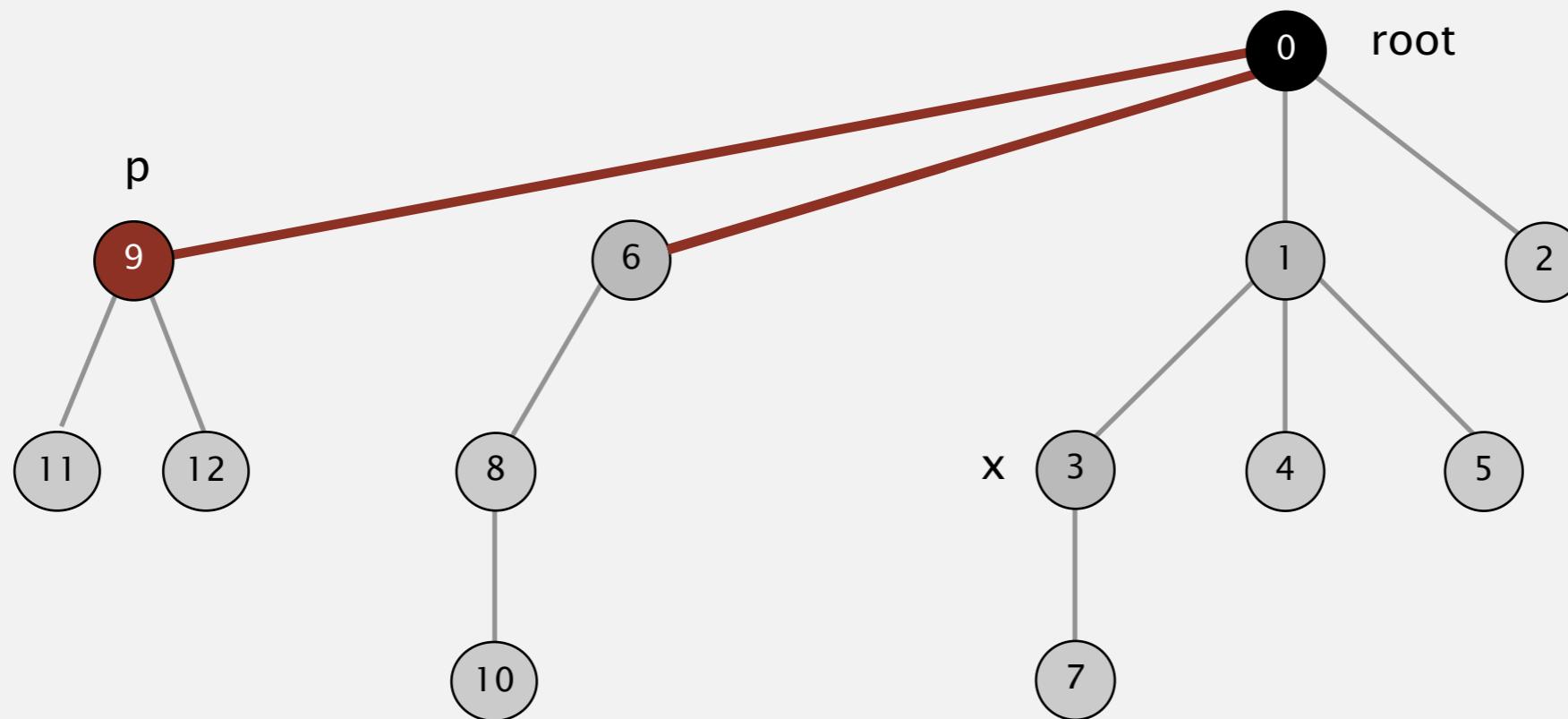
Improvement 2: path compression

Quick union with path compression. Just after computing the root of p , set the `id[]` of each examined node to point to that root.



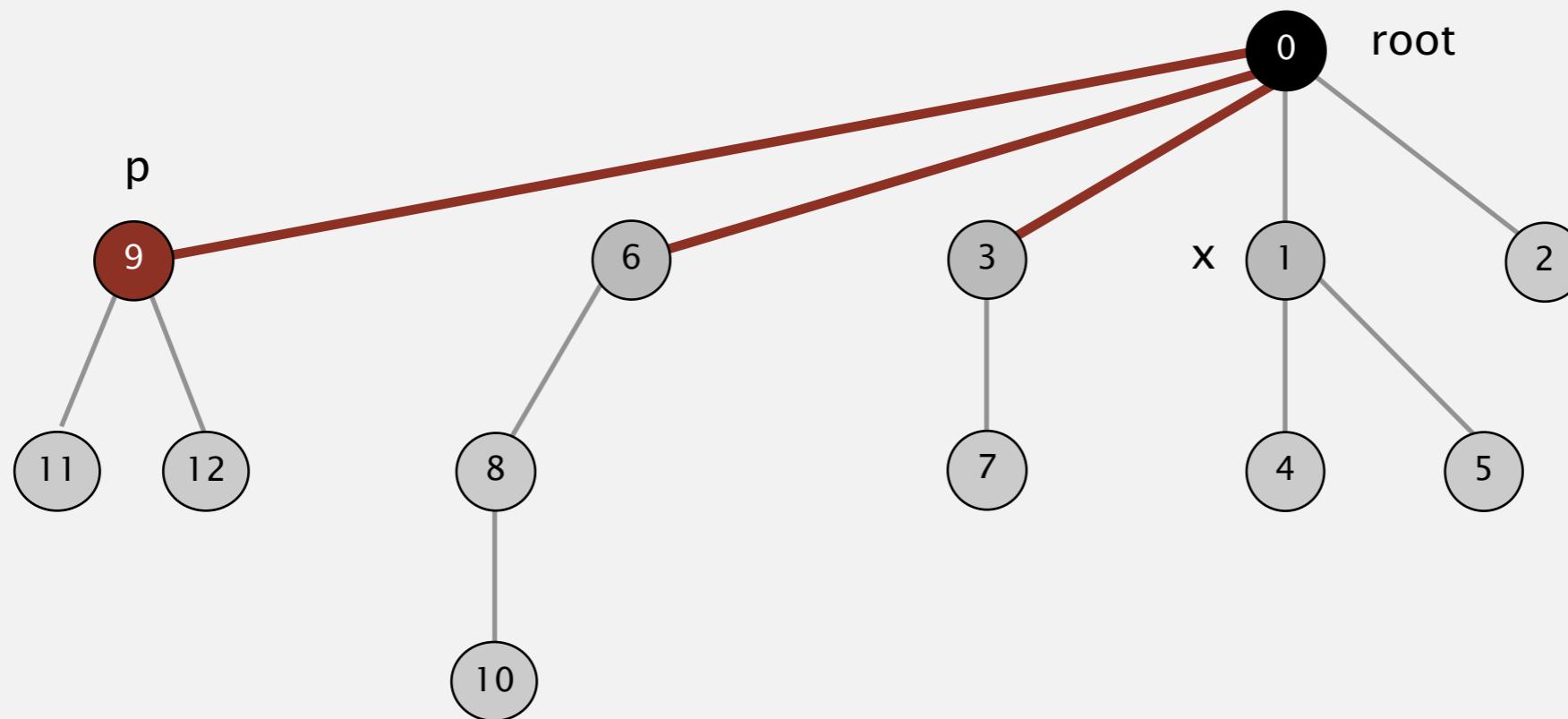
Improvement 2: path compression

Quick union with path compression. Just after computing the root of p , set the `id[]` of each examined node to point to that root.



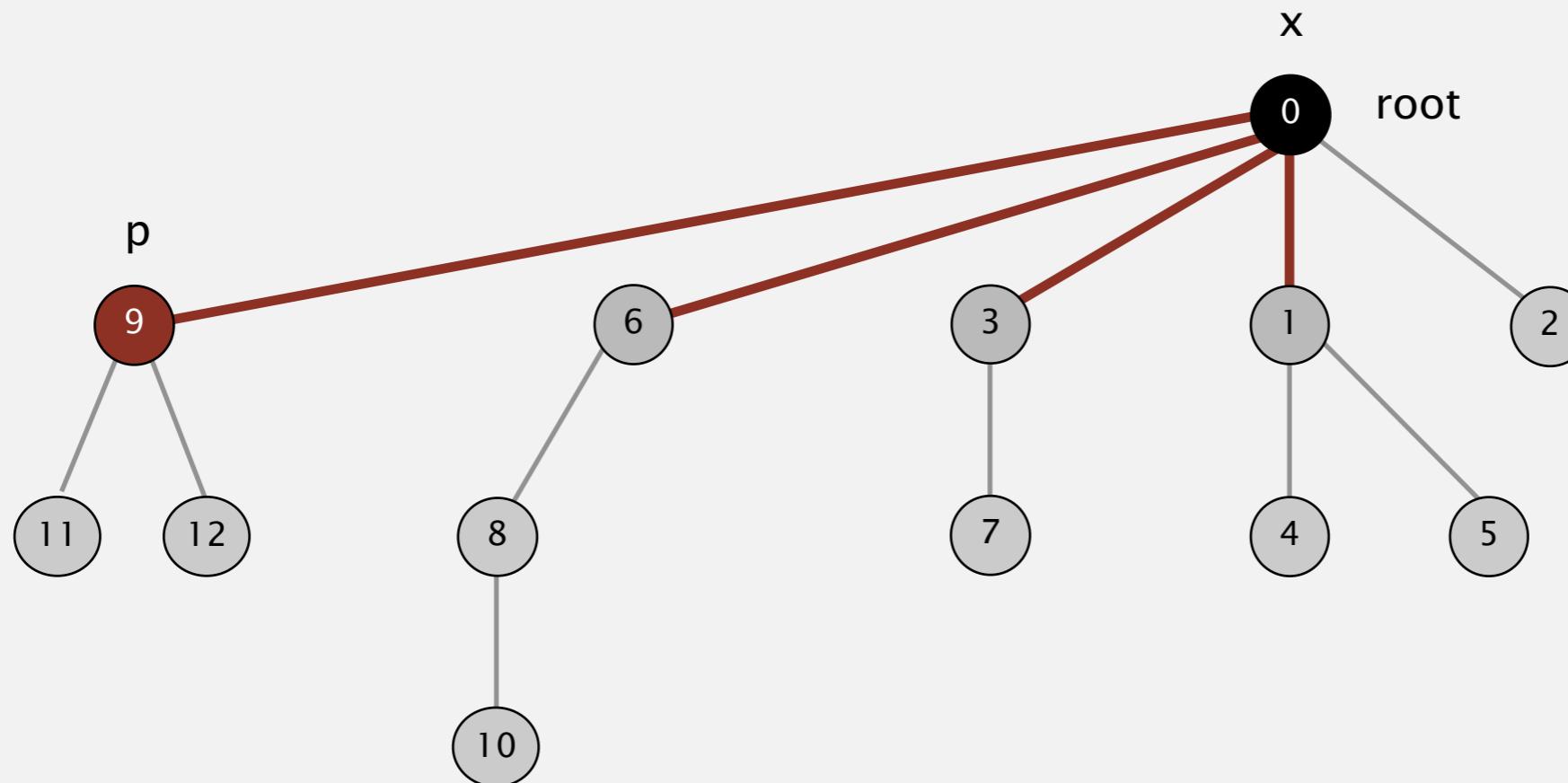
Improvement 2: path compression

Quick union with path compression. Just after computing the root of p , set the `id[]` of each examined node to point to that root.



Improvement 2: path compression

Quick union with path compression. Just after computing the root of p , set the `id[]` of each examined node to point to that root.



Bottom line. Now, `find()` has the side effect of compressing the tree.

Path compression: Java implementation

Two-pass implementation: add second loop to find() to set the id[] of each examined node to the root.

Simpler one-pass variant (path halving): Make every other node in path point to its grandparent.

```
public int find(int i)
{
    while (i != id[i])
    {
        id[i] = id[id[i]]; ← only one extra line of code !
        i = id[i];
    }
    return i;
}
```

In practice. No reason not to! Keeps tree almost completely flat.

Weighted quick-union with path compression: amortized analysis

Proposition. [Hopcroft-Ulman, Tarjan] Starting from an empty data structure, any sequence of M union–find ops on N objects makes $\leq c(N + M \lg^* N)$ array accesses.

- Analysis can be improved to $N + M \alpha(M, N)$.
- Simple algorithm with fascinating mathematics.

N	$\lg^* N$
1	0
2	1
4	2
16	3
65536	4
2^{65536}	5

iterated lg function

Linear-time algorithm for M union–find ops on N objects?

- Cost within constant factor of reading in the data.
- In theory, WQUPC is not quite linear.
- In practice, WQUPC is linear.

Amazing fact. [Fredman-Saks] No linear-time algorithm exists.

Summary

Key point. Weighted quick union (and/or path compression) makes it possible to solve problems that could not otherwise be addressed.

algorithm	worst-case time
quick-find	$M N$
quick-union	$M N$
weighted QU	$N + M \log N$
QU + path compression	$N + M \log N$
weighted QU + path compression	$N + M \lg^* N$

order of growth for M union-find operations on a set of N objects

Ex. [10^9 unions and finds with 10^9 objects]

- WQUPC reduces time from 30 years to 6 seconds.
- Supercomputer won't help much; good algorithm enables solution.

Algorithms

ROBERT SEDGEWICK | KEVIN WAYNE

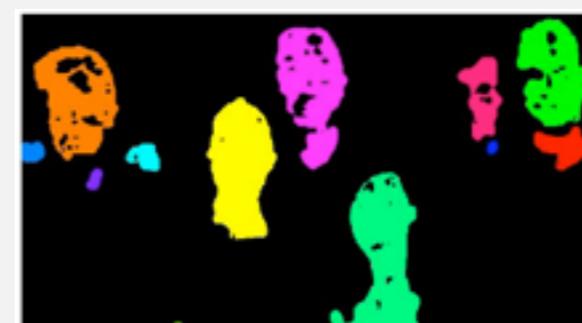
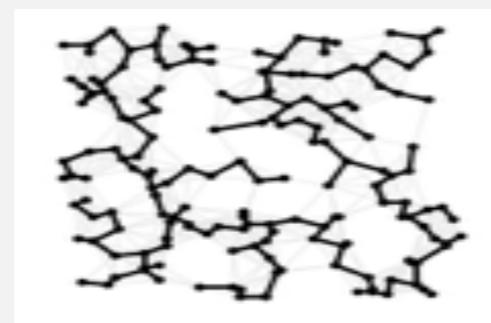
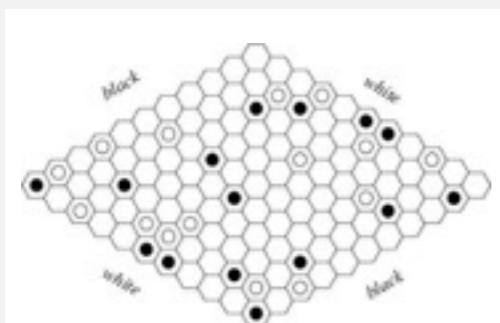
<http://algs4.cs.princeton.edu>

1.5 UNION-FIND

- ▶ *dynamic connectivity*
- ▶ *quick find*
- ▶ *quick union*
- ▶ *improvements*
- ▶ ***applications***

Union-find applications

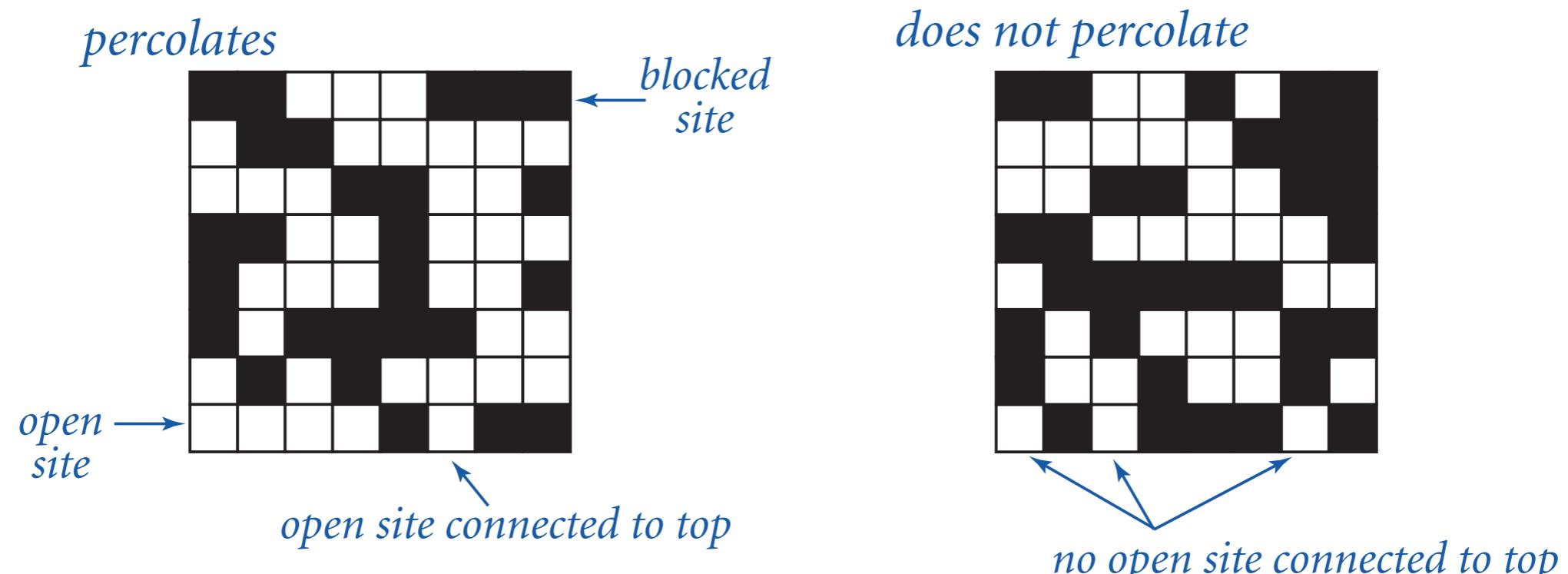
- Percolation.
- Games (Go, Hex).
- ✓ Dynamic connectivity.
 - Least common ancestor.
 - Equivalence of finite state automata.
 - Hoshen-Kopelman algorithm in physics.
 - Hinley-Milner polymorphic type inference.
 - Kruskal's minimum spanning tree algorithm.
 - Compiling equivalence statements in Fortran.
 - Morphological attribute openings and closings.
 - Matlab's bwlabel() function in image processing.



Percolation

An abstract model for many physical systems:

- N -by- N grid of sites.
- Each site is open with probability p (and blocked with probability $1 - p$).
- System **percolates** iff top and bottom are connected by open sites.



Percolation

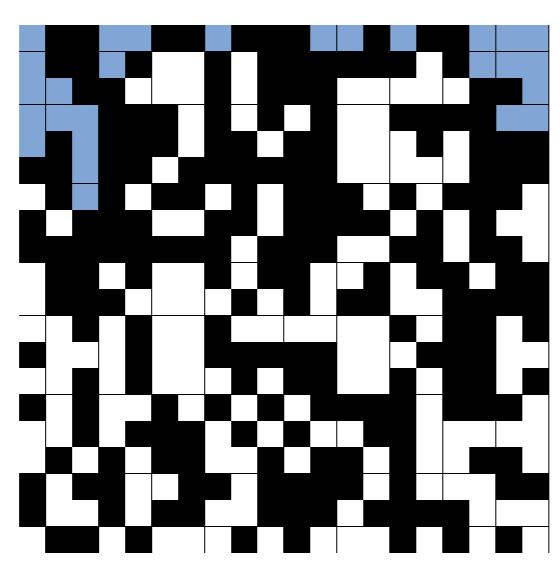
An abstract model for many physical systems:

- N -by- N grid of sites.
- Each site is open with probability p (and blocked with probability $1 - p$).
- System **percolates** iff top and bottom are connected by open sites.

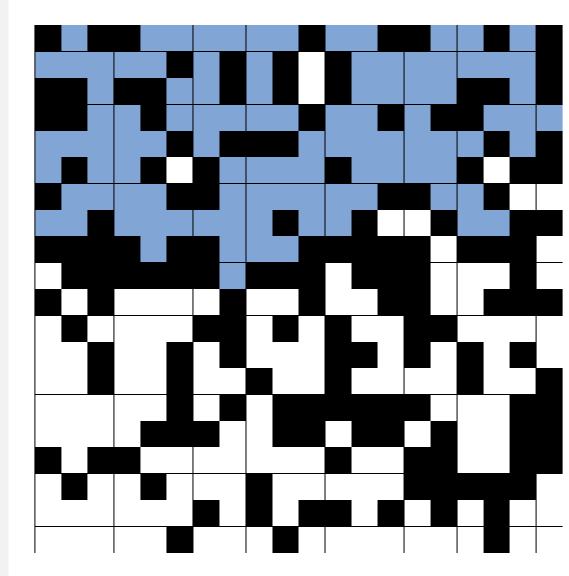
model	system	vacant site	occupied site	percolates
electricity	material	conductor	insulated	conducts
fluid flow	material	empty	blocked	porous
social interaction	population	person	empty	communicates

Likelihood of percolation

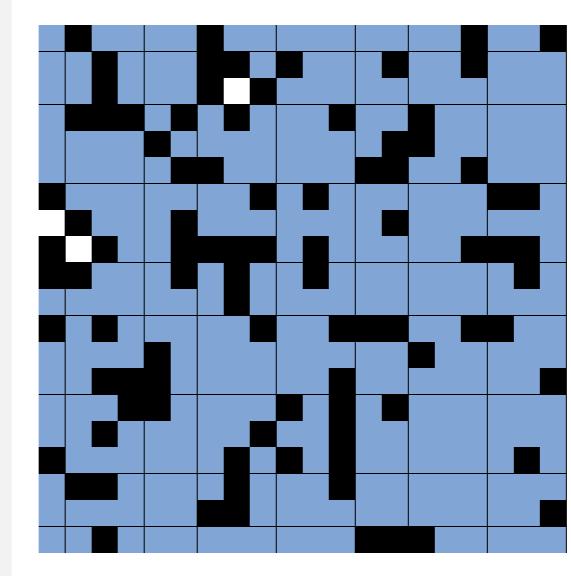
Depends on grid size N and site vacancy probability p .



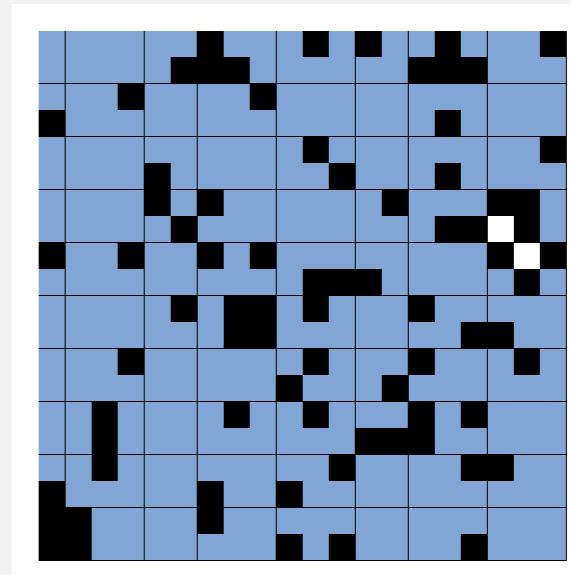
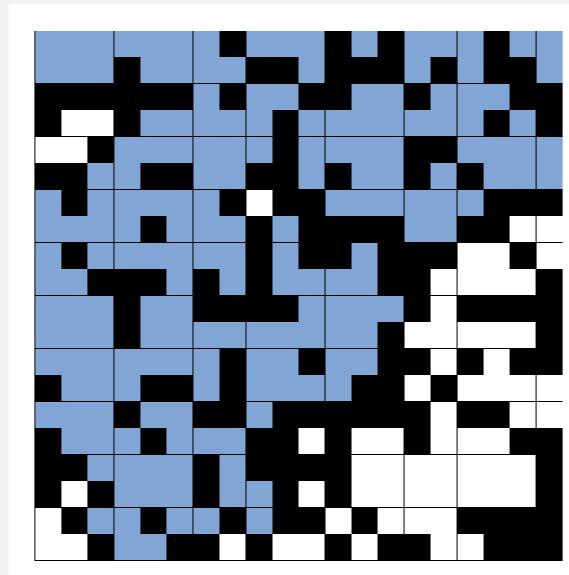
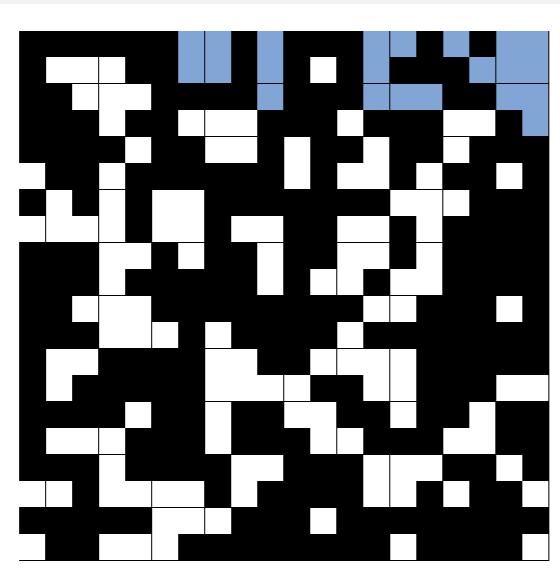
p low (0.4)
does not percolate



p medium (0.6)
percolates?



p high (0.8)
percolates

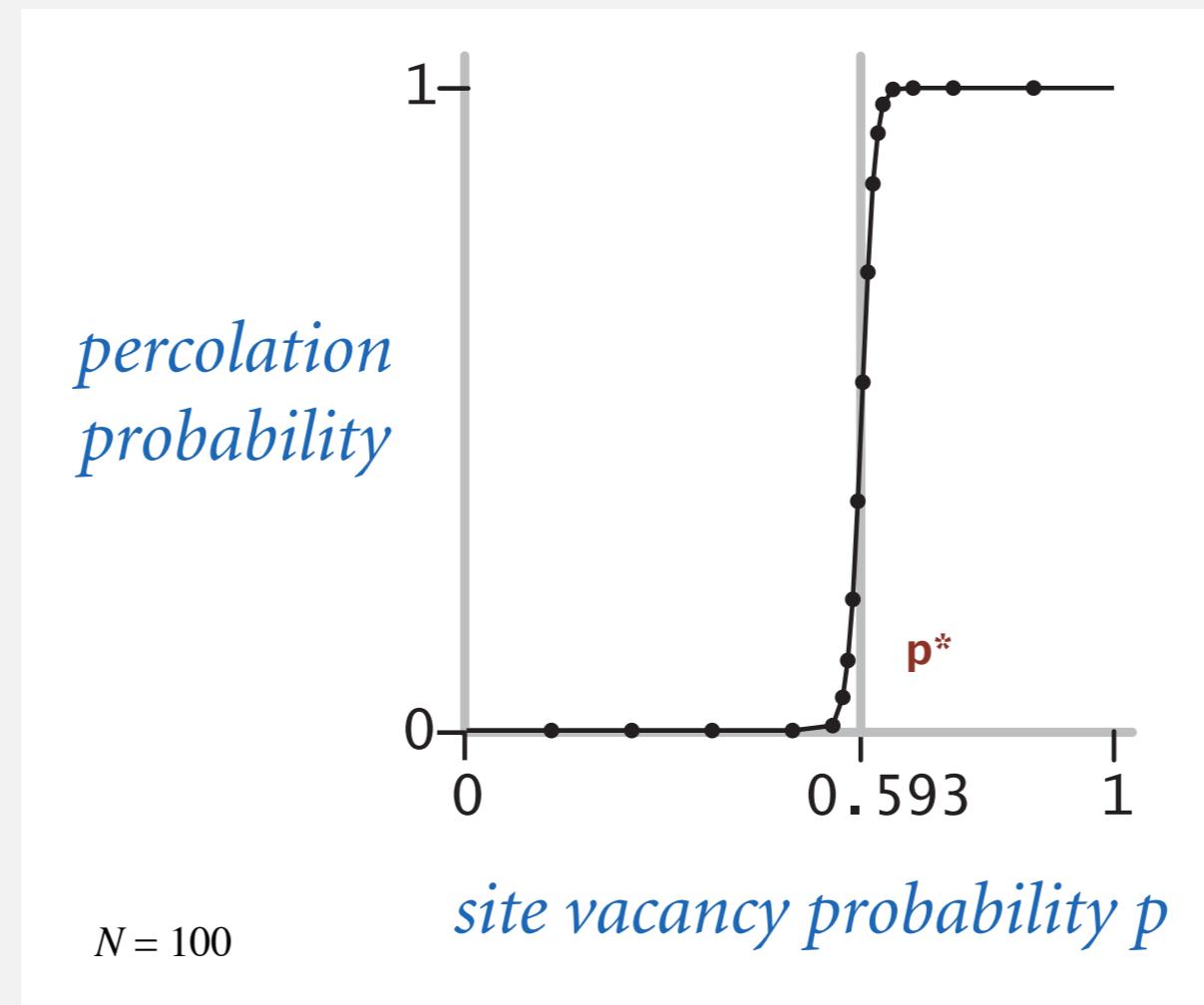


Percolation phase transition

When N is large, theory guarantees a sharp threshold p^* .

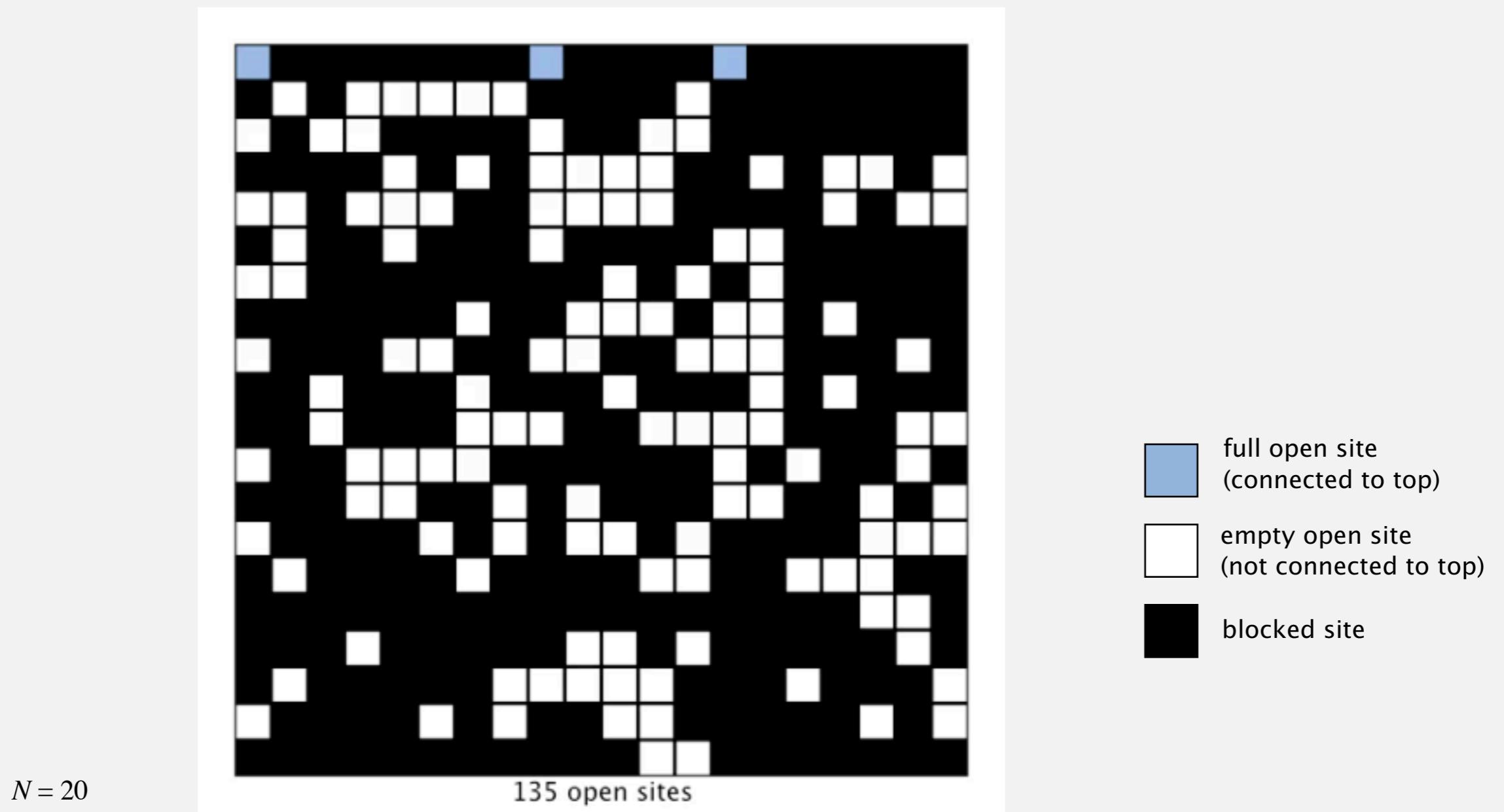
- $p > p^*$: almost certainly percolates.
- $p < p^*$: almost certainly does not percolate.

Q. What is the value of p^* ?



Monte Carlo simulation

- Initialize all sites in an N -by- N grid to be blocked.
- Declare random sites open until top connected to bottom.
- Vacancy percentage estimates p^* .

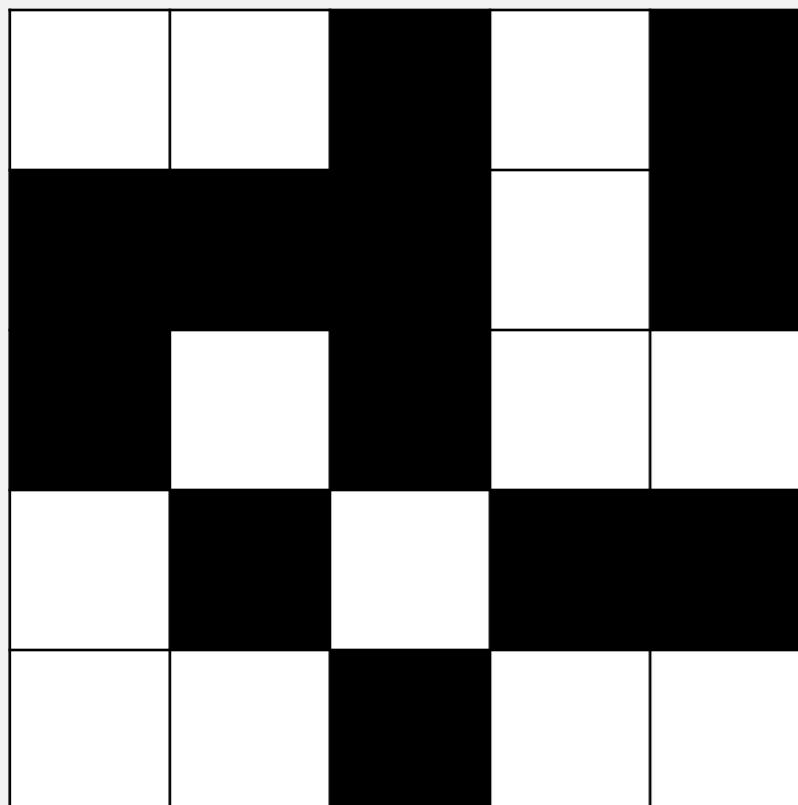


Dynamic connectivity solution to estimate percolation threshold

Q. How to check whether an N -by- N system percolates?

A. Model as a **dynamic connectivity** problem and use **union-find**.

$N = 5$



open site

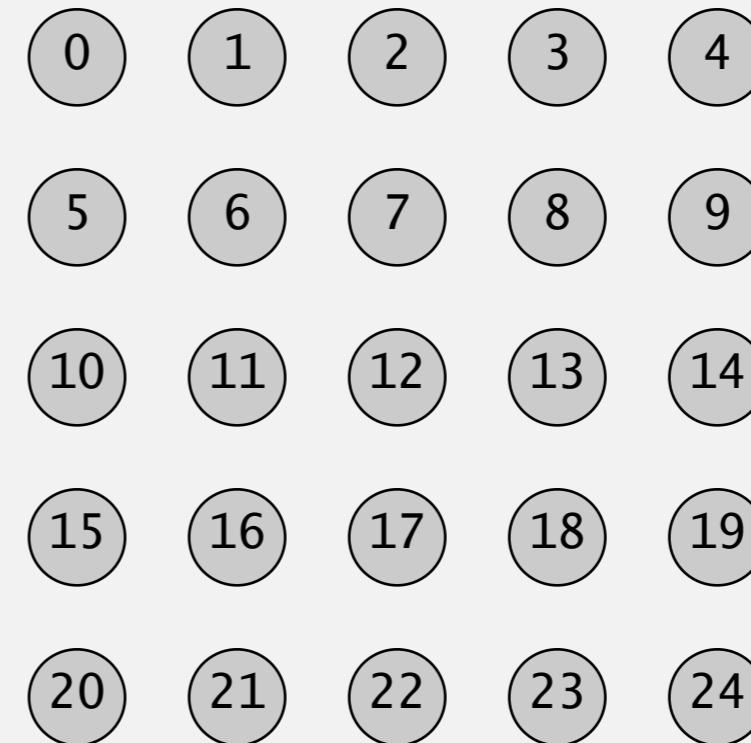
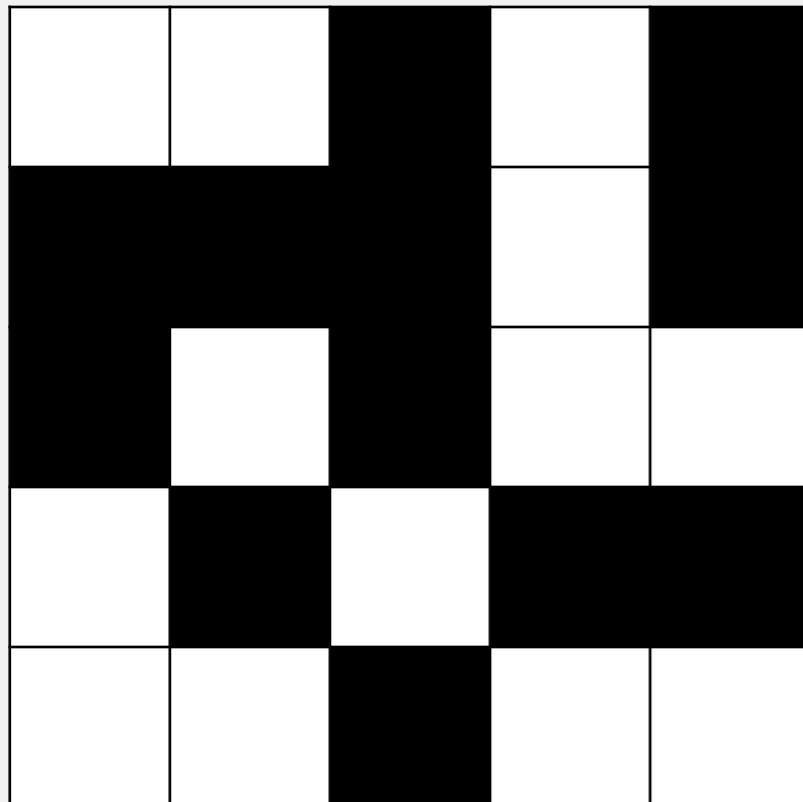
blocked site

Dynamic connectivity solution to estimate percolation threshold

Q. How to check whether an N -by- N system percolates?

- Create an object for each site and name them 0 to $N^2 - 1$.

$N = 5$



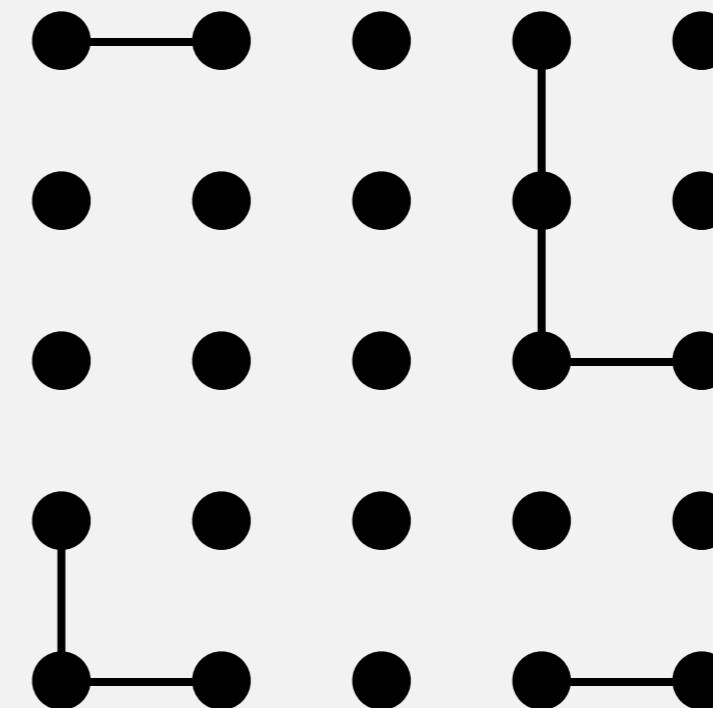
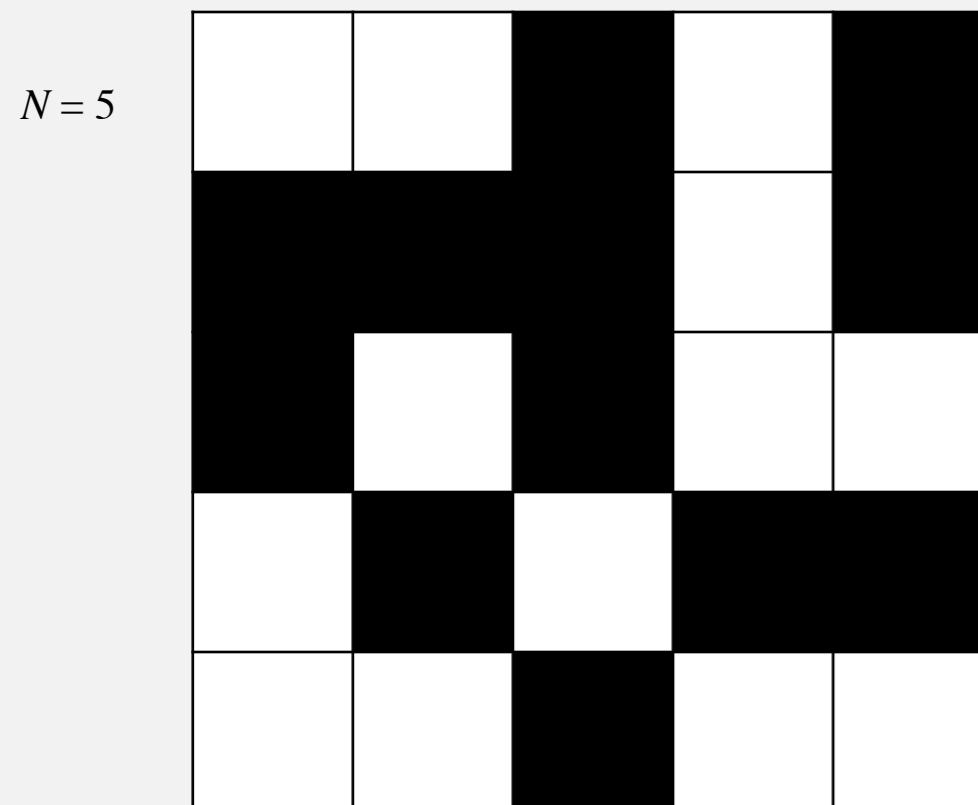
open site

blocked site

Dynamic connectivity solution to estimate percolation threshold

Q. How to check whether an N -by- N system percolates?

- Create an object for each site and name them 0 to $N^2 - 1$.
- Sites are in same component iff connected by open sites.



open site



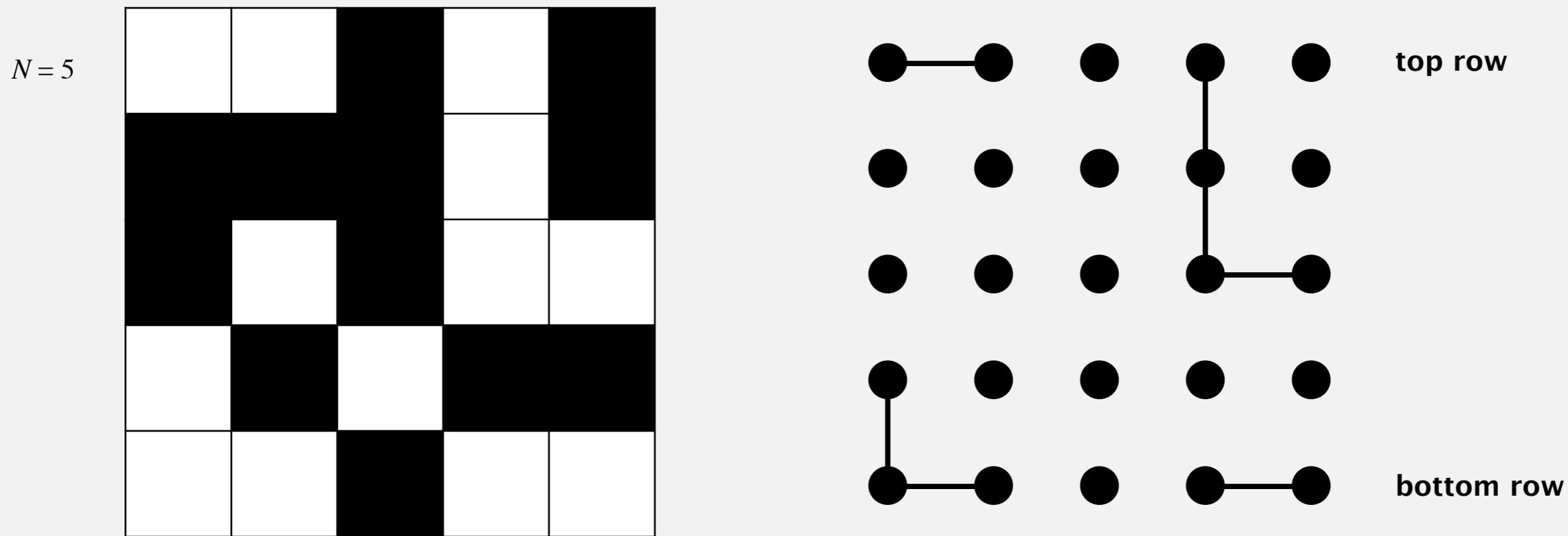
blocked site

Dynamic connectivity solution to estimate percolation threshold

Q. How to check whether an N -by- N system percolates?

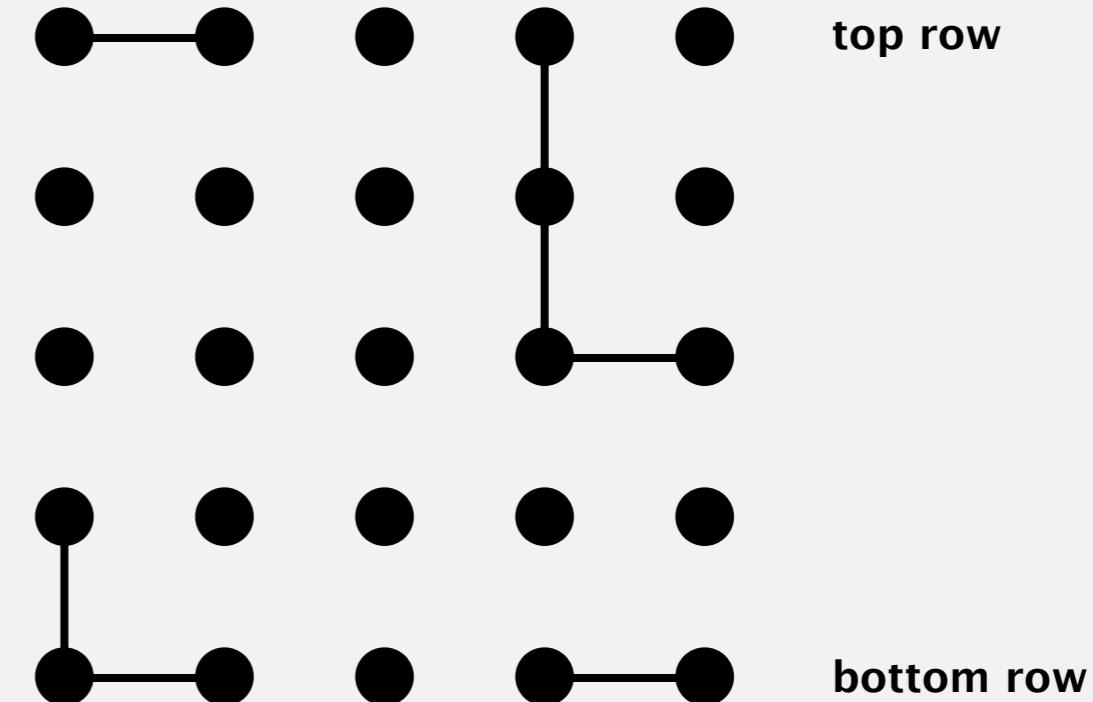
- Create an object for each site and name them 0 to $N^2 - 1$.
- Sites are in same component iff connected by open sites.
- Percolates iff any site on bottom row is connected to any site on top row.

brute-force algorithm: N^2 calls to `connected()`



 open site

 blocked site

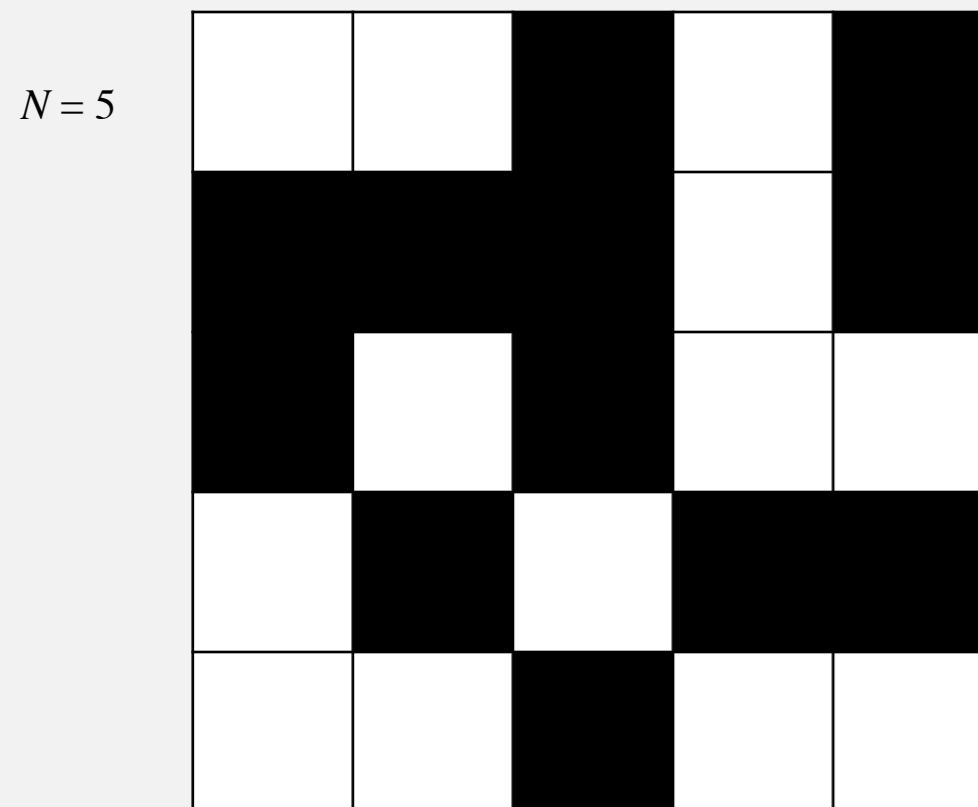


Dynamic connectivity solution to estimate percolation threshold

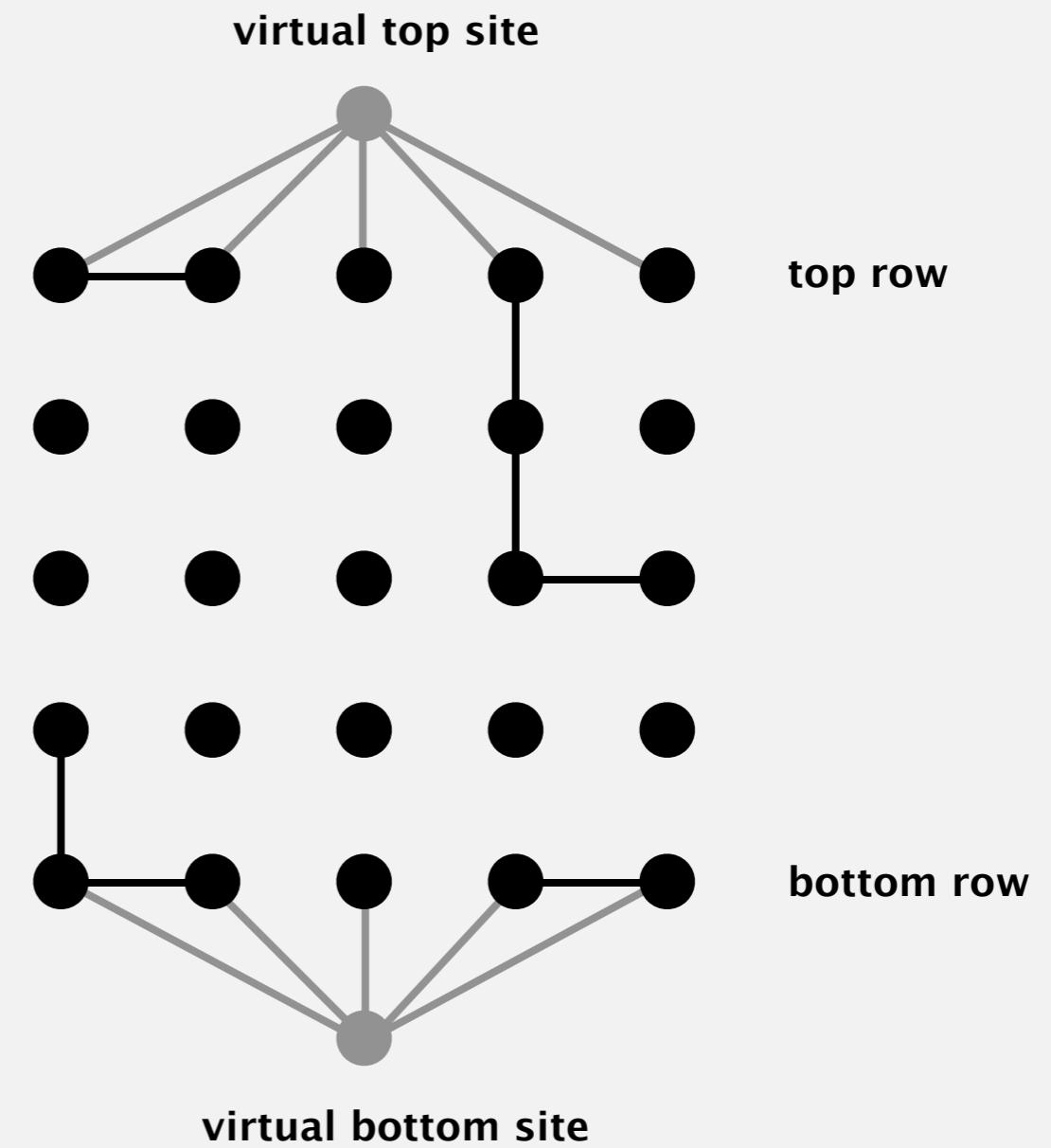
Clever trick. Introduce 2 virtual sites (and connections to top and bottom).

- Percolates iff virtual top site is connected to virtual bottom site.

more efficient algorithm: only 1 call to connected()

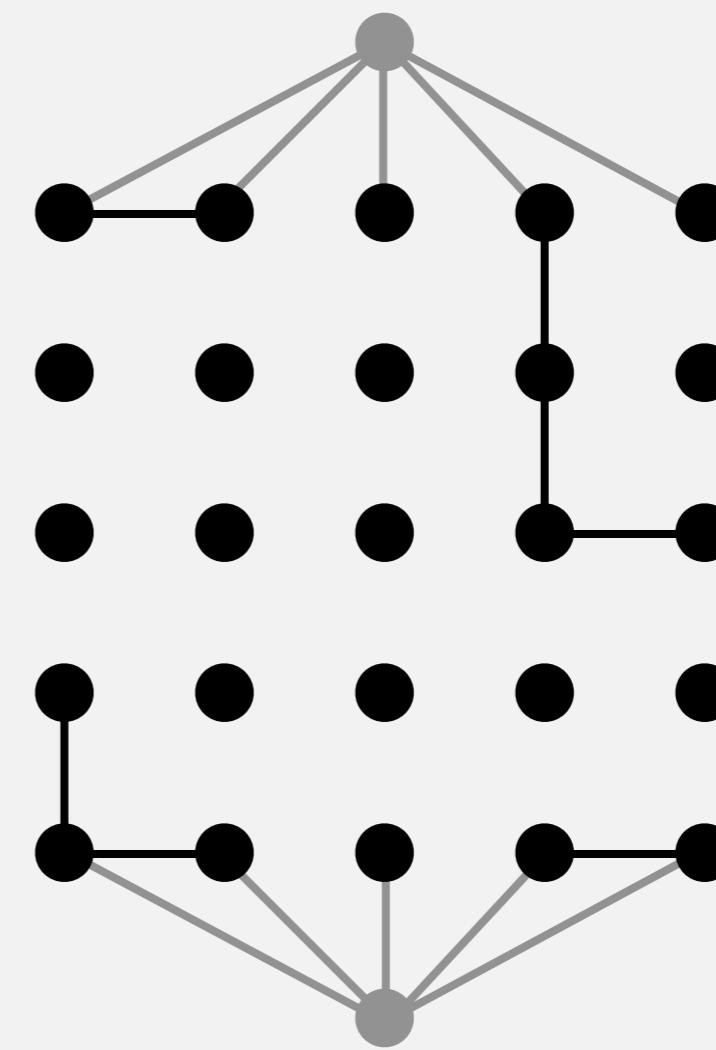
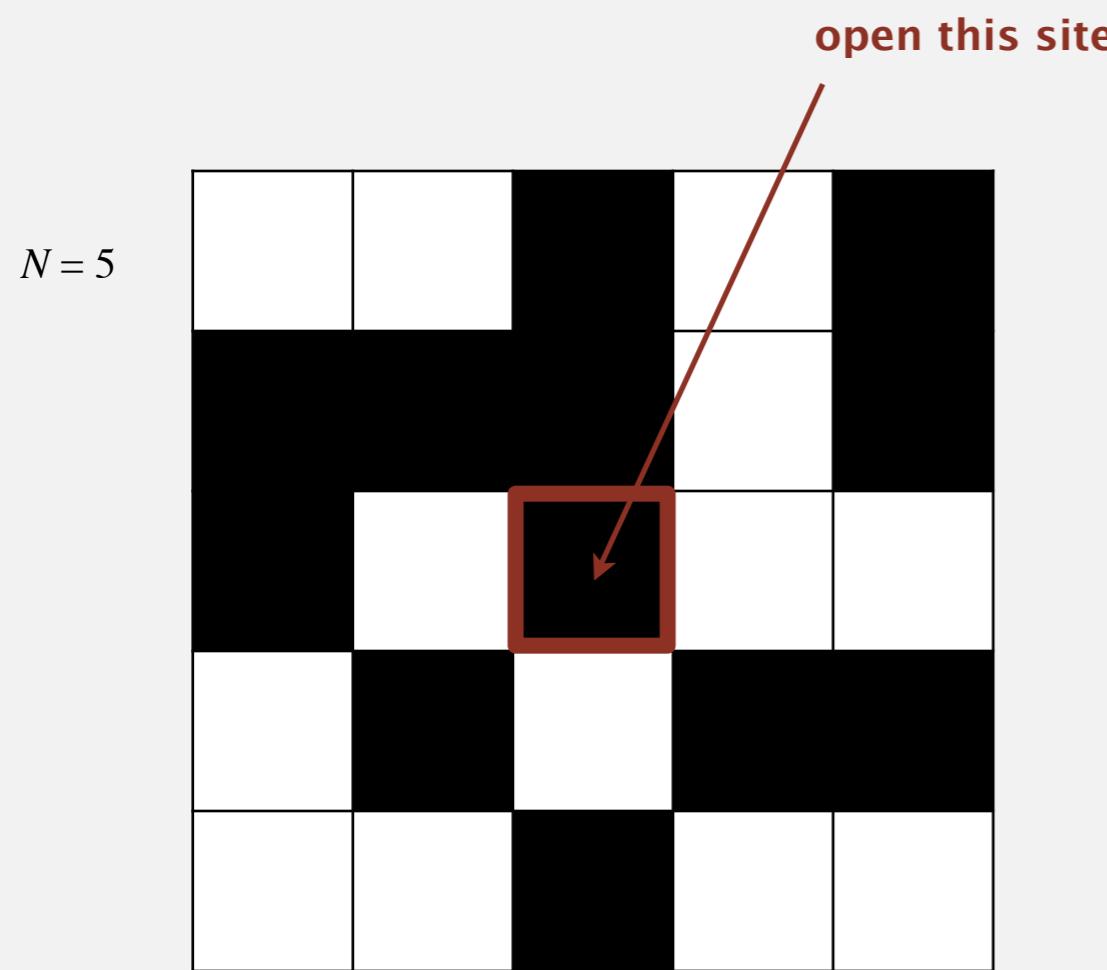


open site
 blocked site



Dynamic connectivity solution to estimate percolation threshold

Q. How to model opening a new site?



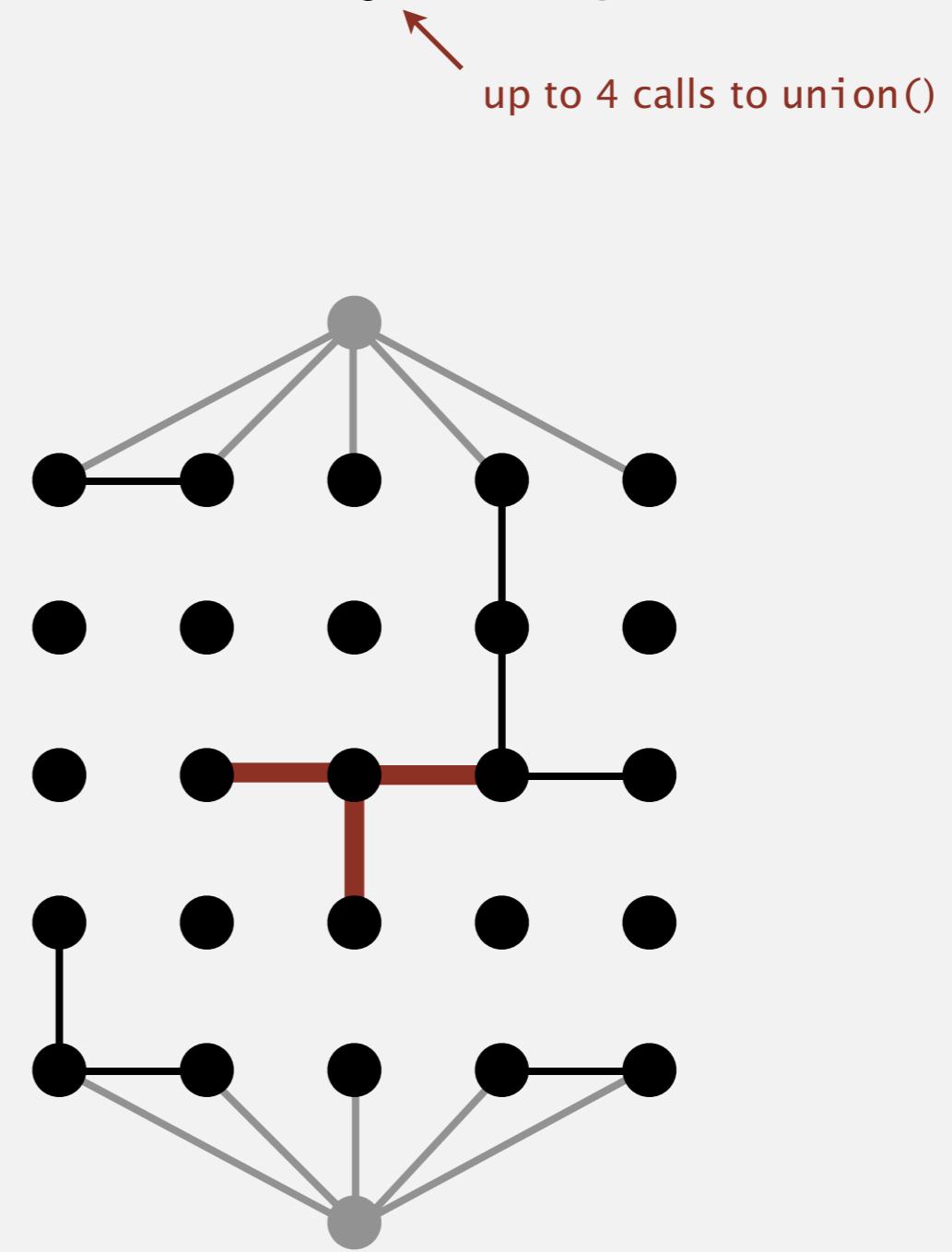
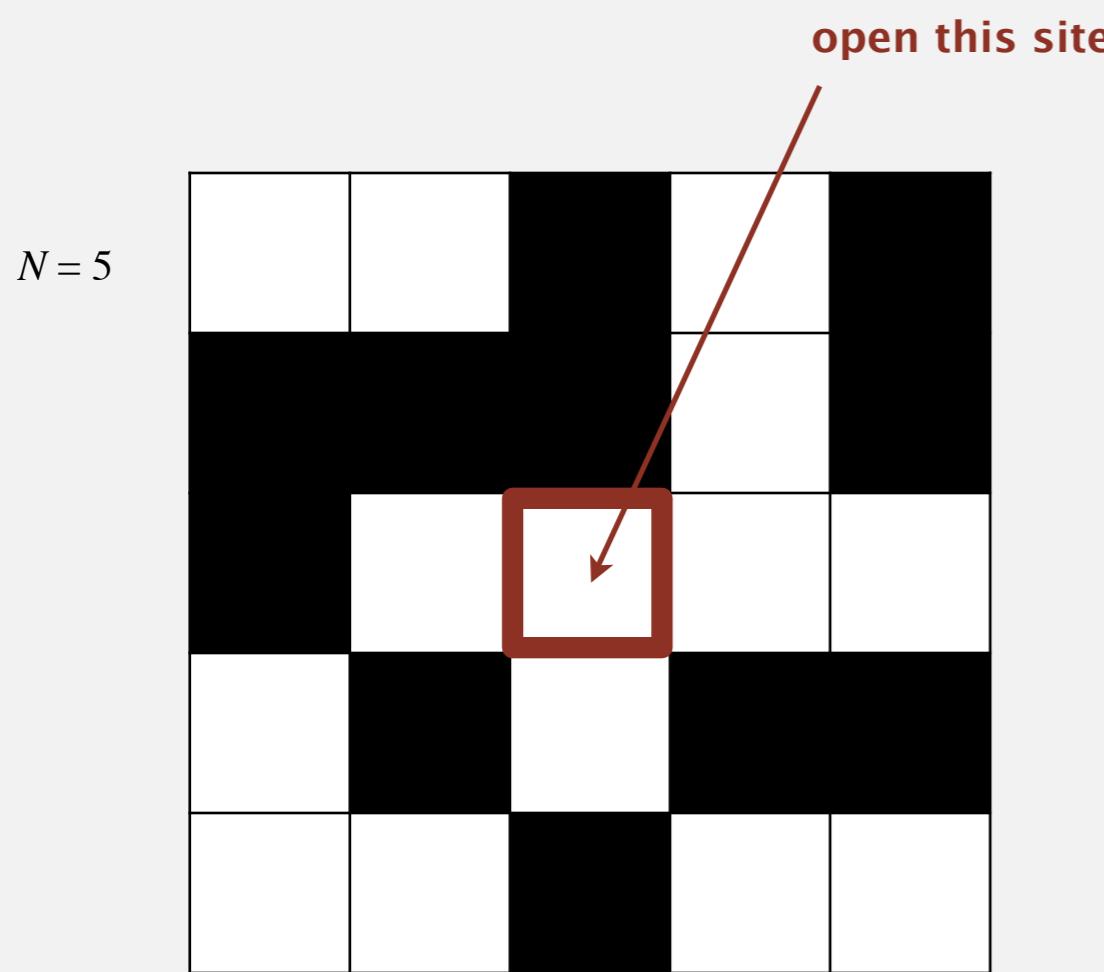
open site

blocked site

Dynamic connectivity solution to estimate percolation threshold

Q. How to model opening a new site?

A. Mark new site as open; connect it to all of its adjacent open sites.



open site

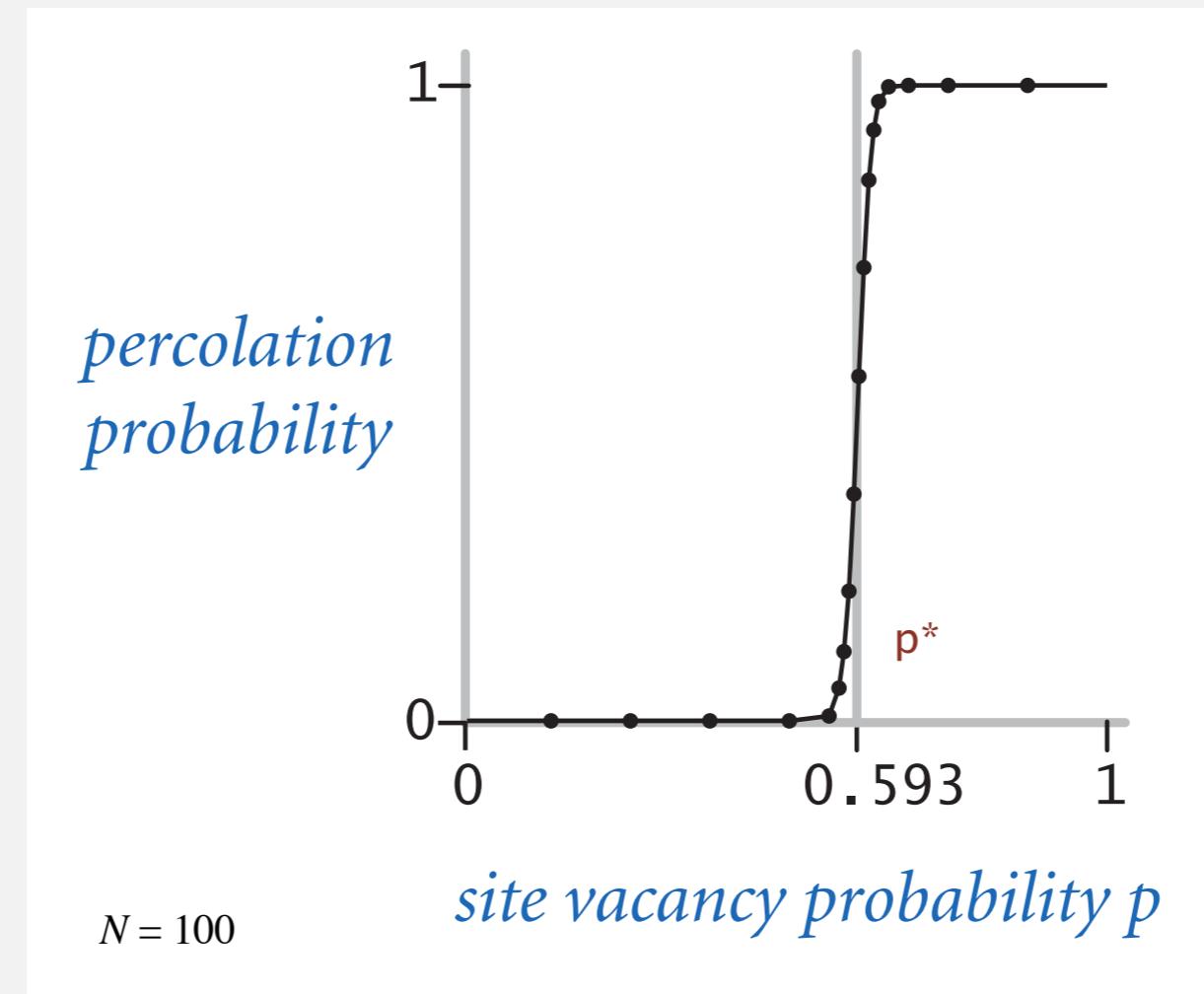
blocked site

Percolation threshold

Q. What is percolation threshold p^* ?

A. About 0.592746 for large square lattices.

constant known only via simulation



Fast algorithm enables accurate answer to scientific question.

Theme of today's lecture (and this course)

Steps to developing a usable algorithm.

- Model the problem.
- Find an algorithm to solve it.
- Fast enough? Fits in memory?
- If not, figure out why.
- Find a way to address the problem.
- Iterate until satisfied.

The scientific method.

Mathematical analysis.

Homework for this week

Buy the textbook

Read the textbook

- Chapter 1.5 covers today's topic
- If you are rusty in Java, please also review Chapters 1.1 and 1.2

Go to recitations

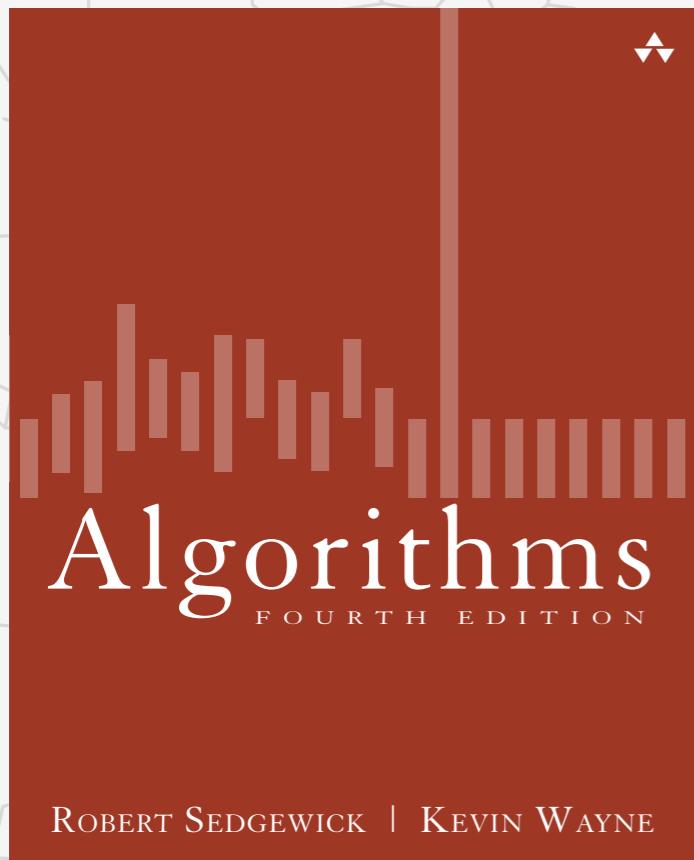
- Recitations start on Monday

Reminders

- Next week's lectures will be by Rajiv Gandhi, reviewing CIS 160
- Your first homework assignment will be released on Tuesday (you can find it on the course homepage)
- <http://www.seas.upenn.edu/~cis121/>

Algorithms

ROBERT SEDGEWICK | KEVIN WAYNE



ROBERT SEDGEWICK | KEVIN WAYNE

<http://algs4.cs.princeton.edu>

1.5 UNION-FIND

- ▶ *dynamic connectivity*
- ▶ *quick find*
- ▶ *quick union*
- ▶ *improvements*
- ▶ *applications*