

COS 226, Fall 2010

*Algorithms
and
Data Structures*



Kevin Wayne

Course Overview

- ▶ **outline**
- ▶ **why study algorithms?**
- ▶ **usual suspects**
- ▶ **coursework**
- ▶ **resources**

COS 226 course overview

What is COS 226?

- Intermediate-level survey course.
- 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, union-find, priority queue
sorting	quicksort, mergesort, heapsort, radix sorts
searching	hash table, BST, red-black tree
graphs	BFS, DFS, Prim, Kruskal, Dijkstra
strings	KMP, regular expressions, TST, Huffman, LZW
geometry	Graham scan, k-d tree, Voronoi diagram

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, ...

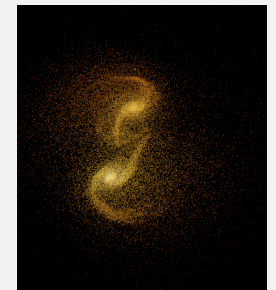
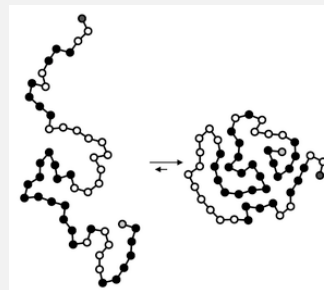
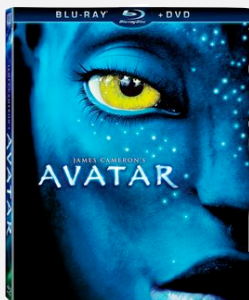
Security. Cell phones, e-commerce, voting machines, ...

Multimedia. CD player, DVD, MP3, JPG, DivX, HDTV, ...

Transportation. Airline crew scheduling, map routing, ...

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

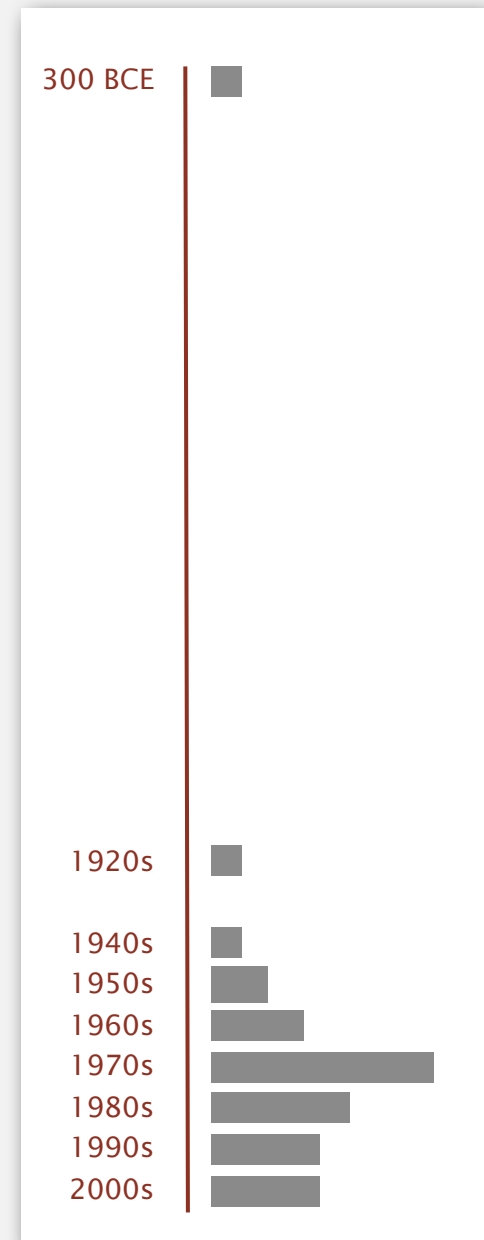
...



Why study algorithms?

Old roots, new opportunities.

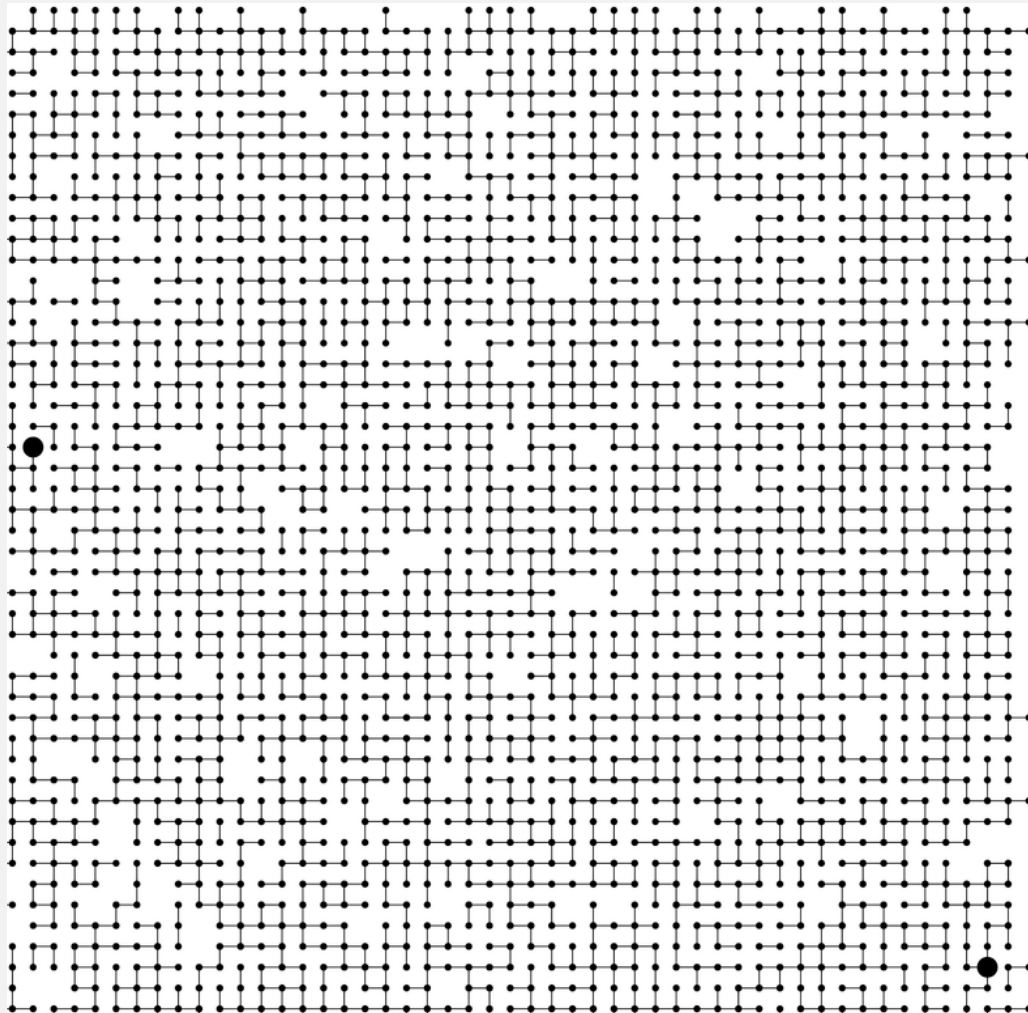
- Study of algorithms dates at least to Euclid.
- Some important algorithms were discovered by undergraduates!



Why study algorithms?

To solve problems that could not otherwise be addressed.

Ex. Network connectivity. [stay tuned]



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. ” — F. Sullivan

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

Why study algorithms?

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

Computational models are replacing mathematical models in scientific inquiry.

$$\begin{aligned} E &= mc^2 \\ F &= ma \end{aligned} \quad F = \frac{Gm_1m_2}{r^2}$$
$$\left[-\frac{\hbar^2}{2m} \nabla^2 + V(r) \right] \Psi(r) = E \Psi(r)$$

20th century science
(formula based)

```
for (double t = 0.0; true; t = t + dt)
  for (int i = 0; i < N; i++)
  {
    bodies[i].resetForce();
    for (int j = 0; j < N; j++)
      if (i != j)
        bodies[i].addForce(bodies[j]);
  }
```

21st century science
(algorithm based)

“Algorithms: a common language for nature, human, and computer.” — A. Wigderson

Why study algorithms?

For fun and profit.

The Google logo, featuring the word "Google" in its signature multi-colored font (blue, red, yellow, blue, green, red) with a trademark symbol.The Nintendo logo, consisting of the word "Nintendo" in a red, rounded font inside a red-outlined rounded rectangle.The Cisco logo, featuring a stylized bridge icon above the word "CISCO" in a red, sans-serif font.The Facebook logo, with the word "facebook" in a white, lowercase, sans-serif font on a dark blue rectangular background.The Adobe logo, featuring a stylized red "A" icon above the word "Adobe" in a black, sans-serif font.The Morgan Stanley logo, with the words "Morgan Stanley" in a white, serif font on a dark blue rectangular background.The RSA logo, featuring the letters "RSA" in a white, bold, sans-serif font on a red rectangular background, with the text "The Security Division of EMC" below it.The Akamai logo, featuring a stylized blue wave icon to the left of the word "Akamai" in an orange, italicized, sans-serif font.The DE Shaw & Co logo, with the text "DE Shaw & Co" in a grey, serif font on a white rectangular background.The Microsoft logo, with the word "Microsoft" in a white, italicized, sans-serif font on a light blue rectangular background.The Yahoo! logo, with the word "YAHOO!" in a red, bold, serif font on a white rectangular background.The Pixar Animation Studios logo, featuring the word "PIXAR" in a large, serif font with a small character holding a lightbulb above the letter "I", and "ANIMATION STUDIOS" in a smaller font below it.

Why study algorithms?

- Their impact is broad and far-reaching.
- Old roots, new opportunities.
- To solve problems that could not otherwise be addressed.
- For intellectual stimulation.
- They may unlock the secrets of life and of the universe.
- For fun and profit.

Why study anything else?



The usual suspects

Lectures. Introduce new material.

Precepts. Discussion, problem-solving, background for programming assignment.

What	When	Where	Who	Office Hours
L01	TTh 11–12:20	Bowen 222	Kevin Wayne	see web
P01	F 11–11:50	Friend 008	Bob Tarjan	see web
P02	F 12:30–1:20	Friend 108	Yuri Pritykin	see web
P02A	F 12:30–1:20	Friend 109	Bob Tarjan	see web
P03	F 1:30–2:20	Friend 008	Aman Dhesi	see web
P03A	F 1:30–2:20	CS 102	Siyu Yang	see web

split into two
(stay tuned)



split into two
(stay tuned)



Computing laboratory. Undergrad TAs in Friend 017. See web for schedule.

Coursework and grading

Programming assignments. 50%

Due at 11pm via electronic submission.

Exams. 20% + 30%

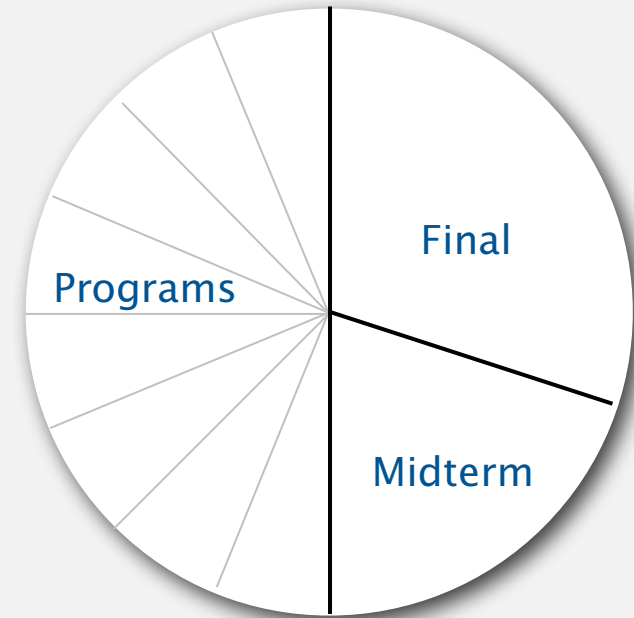
- Closed-book with cheatsheet.
- Midterm (in class on Tuesday, October 26).
- Final (scheduled by Registrar).

Exercises. To be done before precept.

Staff discretion. To adjust borderline cases.

- Precept participation (including discussion of exercises).
- Meet your instructor.
- Report errata.

Check grades? Blackboard.



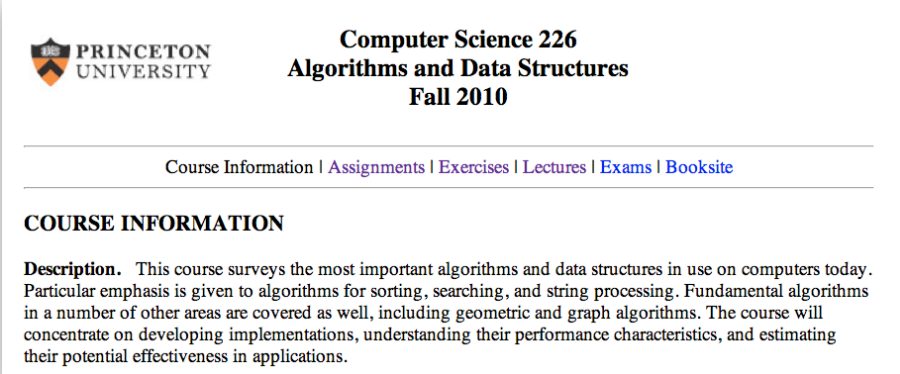
Resources (web)


Course content.

- Course info.
- Programming assignments.
- Exercises.
- Lecture slides.
- Exam archive.
- Submit assignments.

Booksites.

- Brief summary of content.
- Download code from lecture.



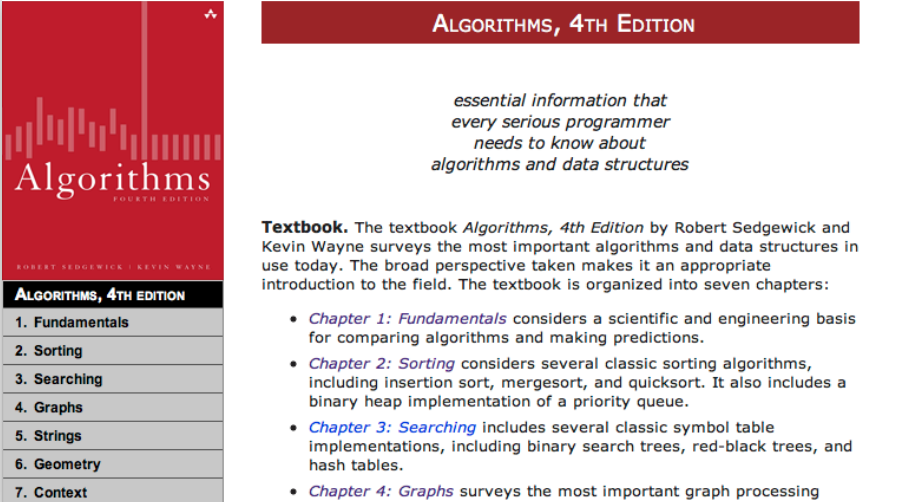
 **Computer Science 226**
Algorithms and Data Structures
Fall 2010

[Course Information](#) | [Assignments](#) | [Exercises](#) | [Lectures](#) | [Exams](#) | [Booksite](#)

COURSE INFORMATION

Description. This course surveys the most important algorithms and data structures in use on computers today. Particular emphasis is given to algorithms for sorting, searching, and string processing. Fundamental algorithms in a number of other areas are covered as well, including geometric and graph algorithms. The course will concentrate on developing implementations, understanding their performance characteristics, and estimating their potential effectiveness in applications.

<http://www.princeton.edu/~cos226>



ALGORITHMS, 4TH EDITION

essential information that every serious programmer needs to know about algorithms and data structures

Textbook. The textbook *Algorithms, 4th Edition* by Robert Sedgwick and Kevin Wayne surveys the most important algorithms and data structures in use today. The broad perspective taken makes it an appropriate introduction to the field. The textbook is organized into seven chapters:

- **Chapter 1: Fundamentals** considers a scientific and engineering basis for comparing algorithms and making predictions.
- **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** includes several classic symbol table implementations, including binary search trees, red-black trees, and hash tables.
- **Chapter 4: Graphs** surveys the most important graph processing

ALGORITHMS, 4TH EDITION
1. Fundamentals
2. Sorting
3. Searching
4. Graphs
5. Strings
6. Geometry
7. Context

<http://www.cs.princeton.edu/IntroProgramming>

<http://www.cs.princeton.edu/Algorithms4>

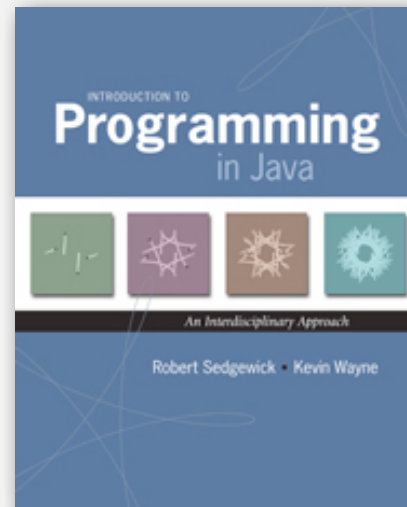
Resources (textbooks)

Required readings.

Algorithms 4th edition, Fall 2010 edition. [only at Labyrinth books]

Recommended Java reference.

Introduction to Programming in Java. [Labyrinth books]



On reserve. Both textbooks are on reserve in the engineering library.

What's ahead?

Lecture 1. Union find. ← today

Precept 1. Meets tomorrow.

Lecture 2. Analysis of algorithms.



Exercise 1. Due in precept tomorrow.

Assignment 1. Due via electronic submission at 11pm on Wednesday.

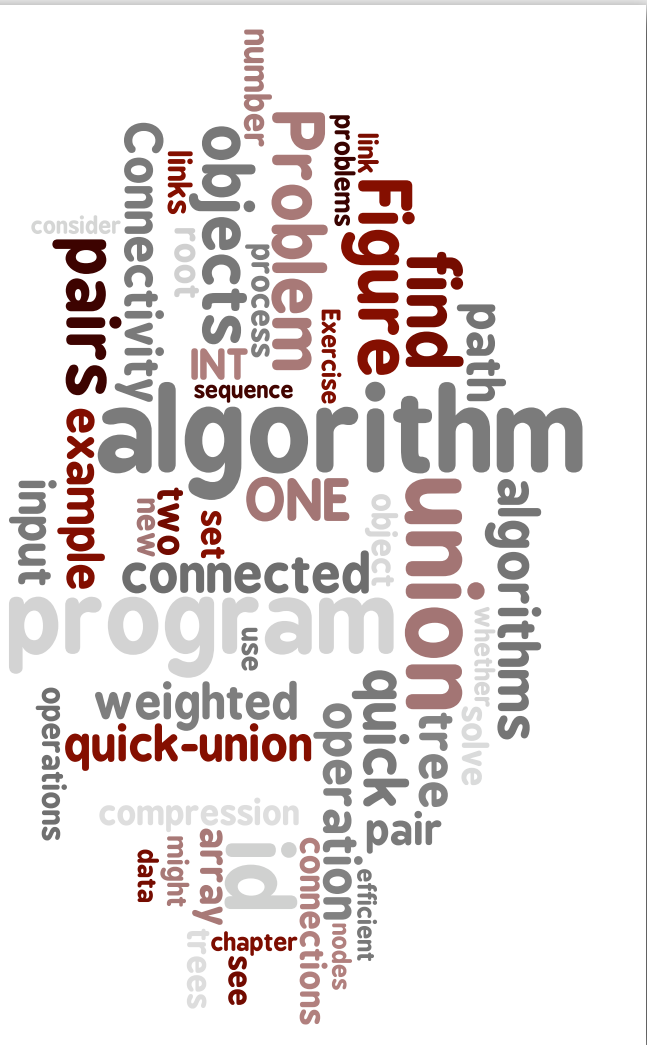
Right course? See me after lecture.

Placed out of COS 126? Review Section 1.5 of Intro to Programming in Java for input-output libraries and command-line interface.

Not registered? Go to any precept tomorrow.

Change precept? Use SCORE. ← see Colleen Kenny-McGinley in CS 210 if the only precept you can attend is closed

1.5 Union Find



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

Subtext 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.

▶ **dynamic connectivity**

- ▶ quick find
- ▶ quick union
- ▶ improvements
- ▶ applications

Dynamic connectivity

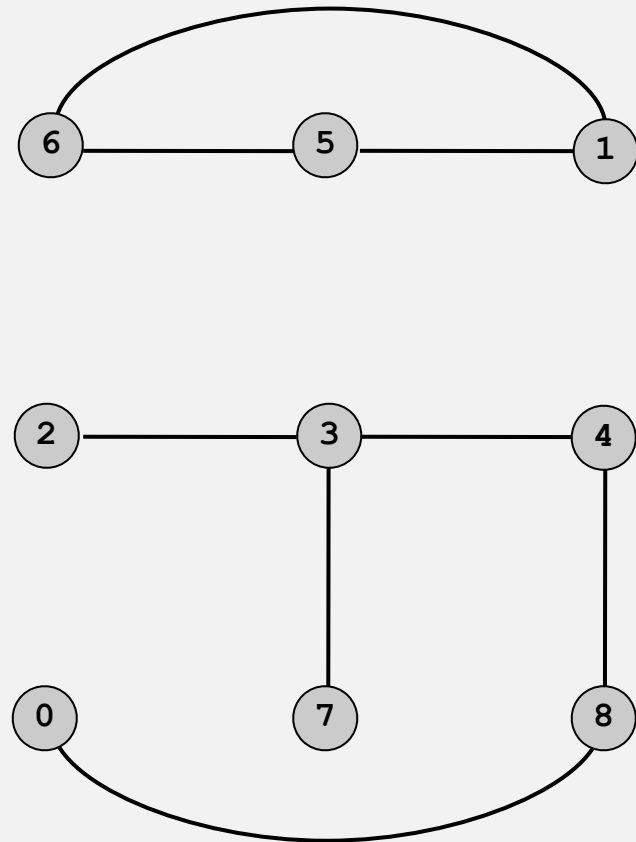
Given a set of objects

- **Union:** connect two objects.
- **Find:** is there a path connecting the two objects?

more difficult problem: find the path

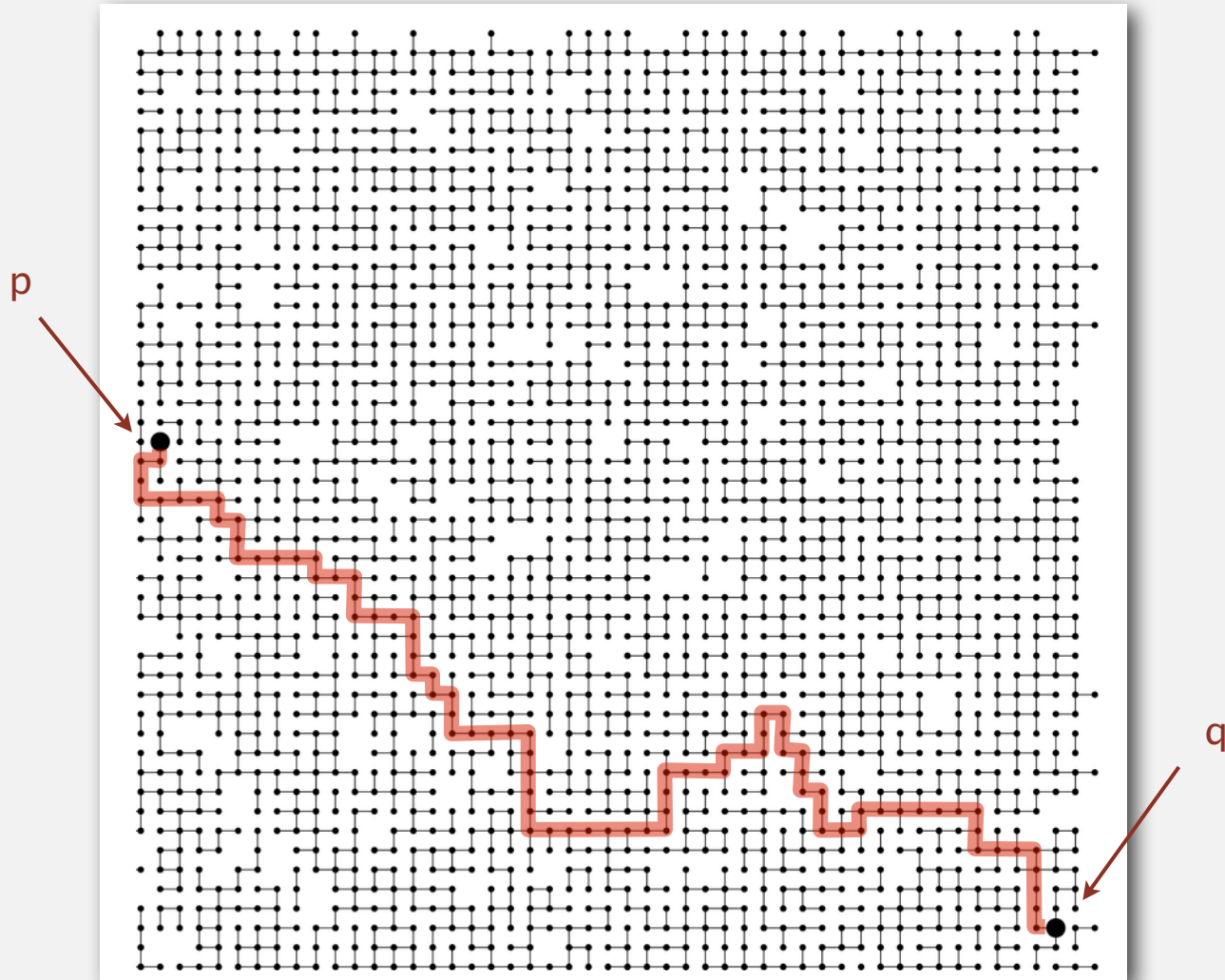


```
union(3, 4)
union(8, 0)
union(2, 3)
union(5, 6)
  find(0, 2)    no
  find(2, 4)    yes
union(5, 1)
union(7, 3)
union(1, 6)
union(4, 8)
  find(0, 2)    yes
  find(2, 4)    yes
```



Connectivity example

Q. Is there a path from p to q ?



A. Yes.

Modeling the objects

Dynamic connectivity applications involve manipulating objects of all types.

- Pixels in a digital photo.
- Computers in a network.
- Variable names in Fortran.
- Friends in a social network.
- Transistors in a computer chip.
- Elements in a mathematical set.
- 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.

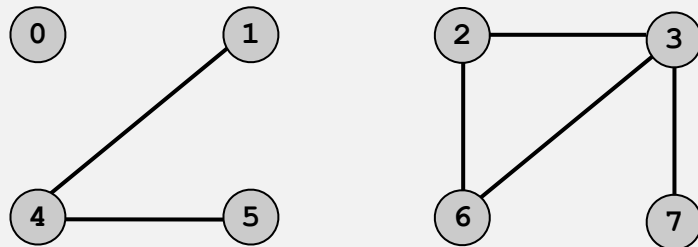
can use symbol table to translate from object names to integers: stay tuned (Chapter 3)

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 components. Maximal **set** of objects that are mutually connected.



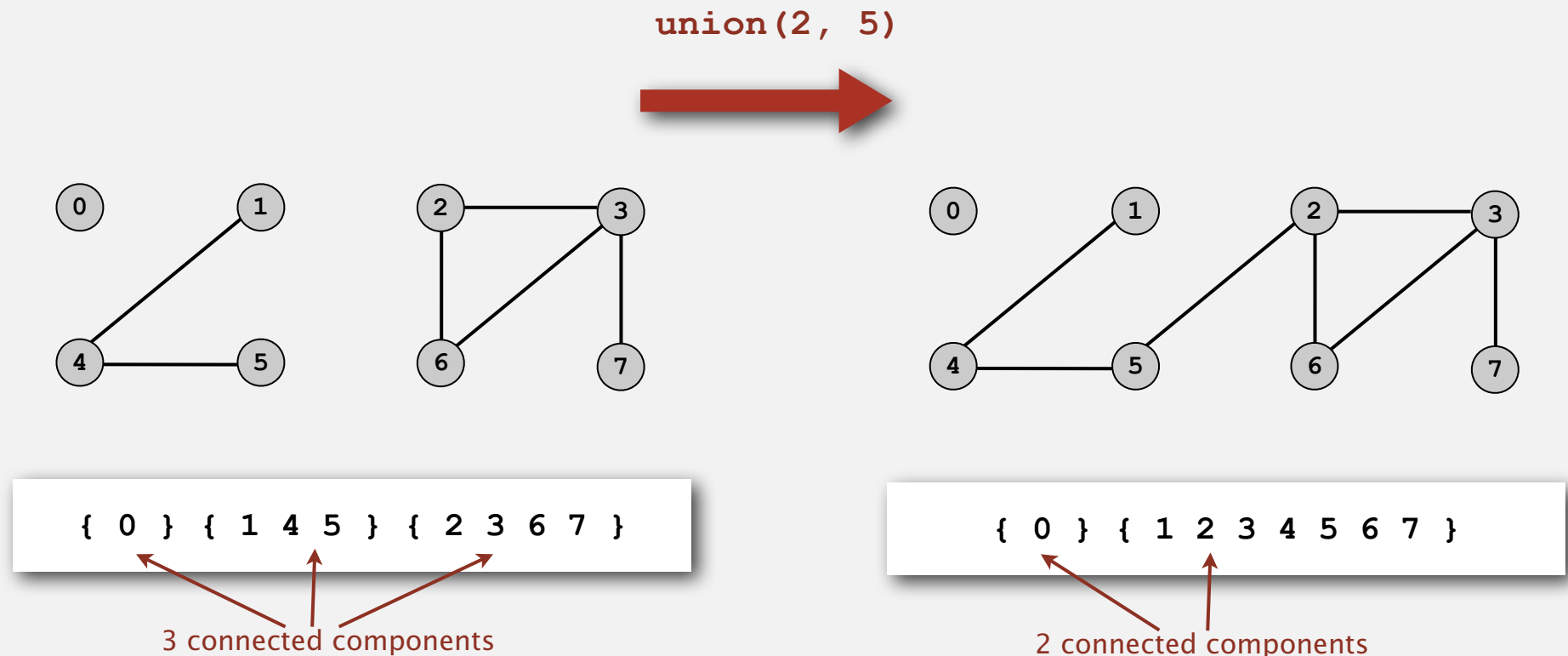
{ 0 } { 1 4 5 } { 2 3 6 7 }

3 connected components

Implementing the operations

Find query. Check if two objects are in the same component.

Union command. Replace components containing two objects with their union.



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.
- Find queries and union commands may be intermixed.

```
public class UF
```

```
    UF(int N)
```

*create union-find data structure with
N objects and no connections*

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

are p and q in the same component?

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

add connection between p and q

```
    int count()
```

number of components

Dynamic connectivity client

- Read in value N .
- Repeat:
 - read in pair of integers
 - write out pair if they are not already connected

```
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.find(p, q)) continue;
        uf.union(p, q);
        StdOut.println(p + " " + q);
    }
}
```

```
% more tiny.txt
10
4 3
3 8
6 5
9 4
2 1
8 9
5 0
7 2
6 1
1 0
6 7
```

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

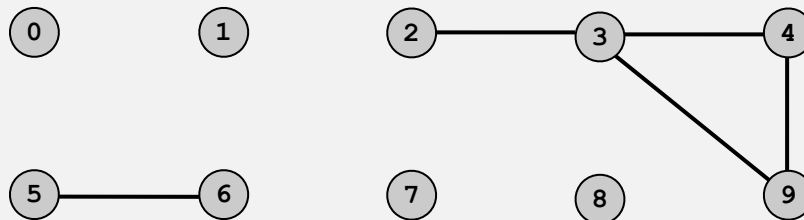
Quick-find [eager approach]

Data structure.

- Integer array `id[]` of size `N`.
- Interpretation: `p` and `q` in same component iff they have the same `id`.

<code>i</code>	0	1	2	3	4	5	6	7	8	9
<code>id[i]</code>	0	1	9	9	9	6	6	7	8	9

5 and 6 are connected
2, 3, 4, and 9 are connected



Quick-find [eager approach]

Data structure.

- Integer array `id[]` of size `N`.
- Interpretation: `p` and `q` in same component iff they have the same `id`.

<code>i</code>	0	1	2	3	4	5	6	7	8	9
<code>id[i]</code>	0	1	9	9	9	6	6	7	8	9

5 and 6 are connected
2, 3, 4, and 9 are connected

Find. Check if `p` and `q` have the same `id`.

`id[3] = 9; id[6] = 6`
3 and 6 in different components

Quick-find [eager approach]

Data structure.

- Integer array $id[]$ of size N .
- Interpretation: p and q in same component iff they have the same id .

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

5 and 6 are connected
2, 3, 4, and 9 are connected

Find. Check if p and q have the same id .

$id[3] = 9; id[6] = 6$
3 and 6 in different components

Union. To merge sets containing p and q , change all entries with $id[p]$ to $id[q]$.

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

problem: many values can change

union of 3 and 6
2, 3, 4, 5, 6, and 9 are connected

Quick-find example

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

id[p] and id[q] differ, so union() changes entries equal to id[p] to id[q] (in red)

id[p] and id[q] match, so no change

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;
```

```
    }
```

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

```
    { return id[p] == id[q]; }
```

```
    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;
```

```
    }
```

```
}
```

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

check whether p and q
are in the same component
(2 array accesses)

change all entries with $id[p]$ to $id[q]$
(N array accesses)

Quick-find is too slow

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

algorithm	init	union	find
quick-find	N	N	1

Quick-find defect.

- Union too expensive.
- Trees are flat, but too expensive to keep them flat.
- Ex. Takes N^2 array accesses to process sequence of N union commands on N objects.

Quadratic algorithms do not scale

Rough standard (for now).

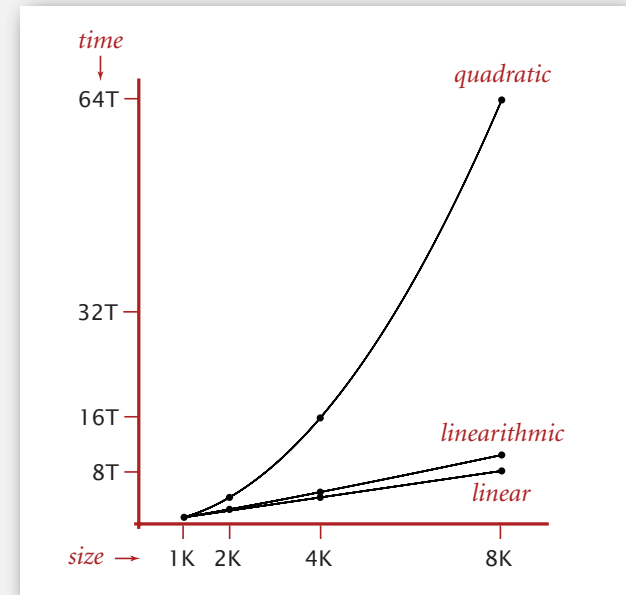
- 10^9 operations per second.
- 10^9 words of main memory.
- Touch all words in approximately 1 second.

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!



Paradoxically, quadratic algorithms get worse with newer equipment.

- New computer may be 10x as fast.
- But, has 10x as much memory so problem may be 10x bigger.
- With quadratic algorithm, takes 10x as long!

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

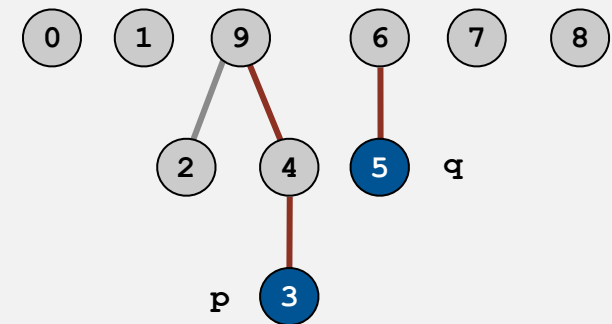
Quick-union [lazy approach]

Data structure.

- Integer array `id[]` of size `N`.
- Interpretation: `id[i]` is parent of `i`.
- **Root** of `i` is `id[id[id[...id[i]...]]]`.

keep going until it doesn't change

<code>i</code>	0	1	2	3	4	5	6	7	8	9
<code>id[i]</code>	0	1	9	4	9	6	6	7	8	9



3's root is 9; 5's root is 6

Quick-union [lazy approach]

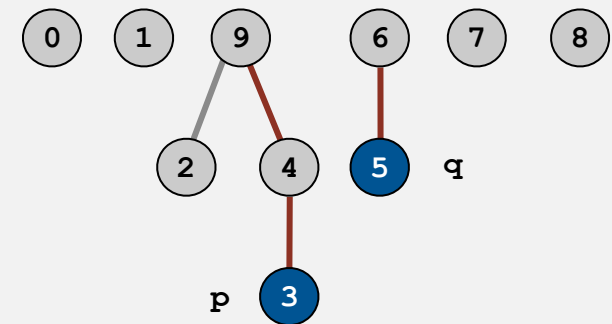
Data structure.

- Integer array `id[]` of size `N`.
- Interpretation: `id[i]` is parent of `i`.
- **Root** of `i` is `id[id[id[...id[i]...]]]`.

keep going until it doesn't change

<code>i</code>	0	1	2	3	4	5	6	7	8	9
<code>id[i]</code>	0	1	9	4	9	6	6	7	8	9

Find. Check if `p` and `q` have the same root.



3's root is 9; 5's root is 6
3 and 5 are in different components

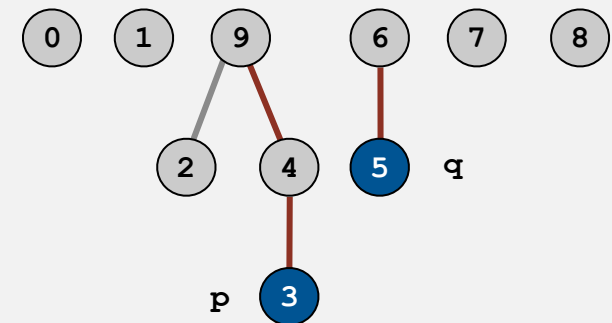
Quick-union [lazy approach]

Data structure.

- Integer array `id[]` of size `N`.
- Interpretation: `id[i]` is parent of `i`.
- **Root** of `i` is `id[id[...id[i]...]]`.

keep going until it doesn't change

<code>i</code>	0	1	2	3	4	5	6	7	8	9
<code>id[i]</code>	0	1	9	4	9	6	6	7	8	9



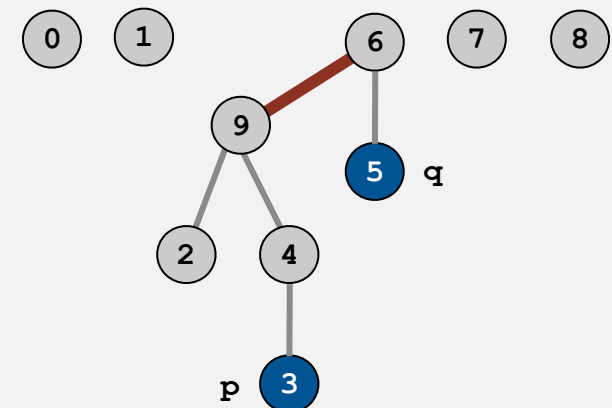
3's root is 9; 5's root is 6
3 and 5 are in different components

Find. Check if `p` and `q` have the same root.

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

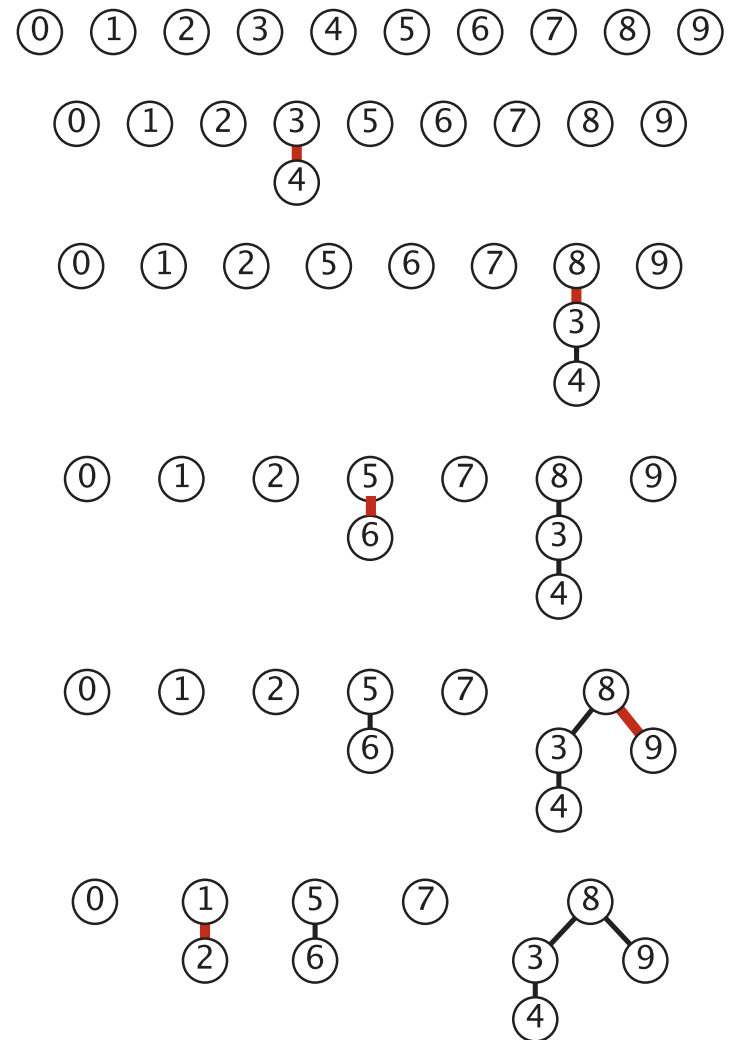
<code>i</code>	0	1	2	3	4	5	6	7	8	9
<code>id[i]</code>	0	1	9	4	9	6	6	7	8	6

only one value changes



Quick-union example

		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



Quick-union example

		id[]									
<u>p</u>	<u>q</u>	0	1	2	3	4	5	6	7	8	9

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

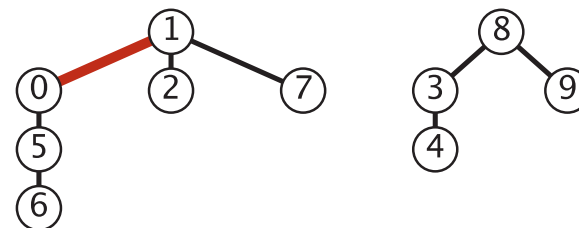
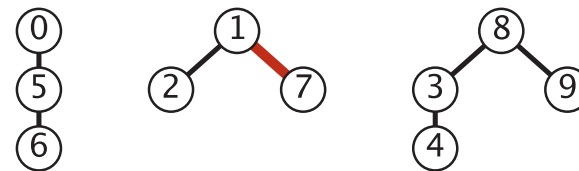
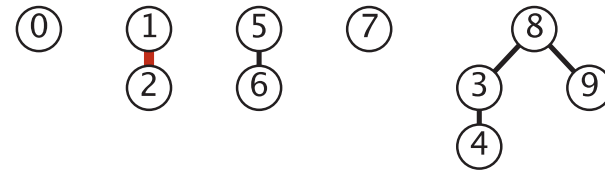
0 1 1 8 3 0 5 1 8 8

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: 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;
    }

    private int root(int i)
    {
        while (i != id[i]) i = id[i];
        return i;
    }

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

    public void union(int p, int q)
    {
        int i = root(p), j = root(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)

check if p and q have same root
(depth of p and q 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	init	union	find
quick-find	N	N	1
quick-union	N	$N \dagger$	N

← worst case

† includes cost of finding root

Quick-find defect.

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

Quick-union defect.

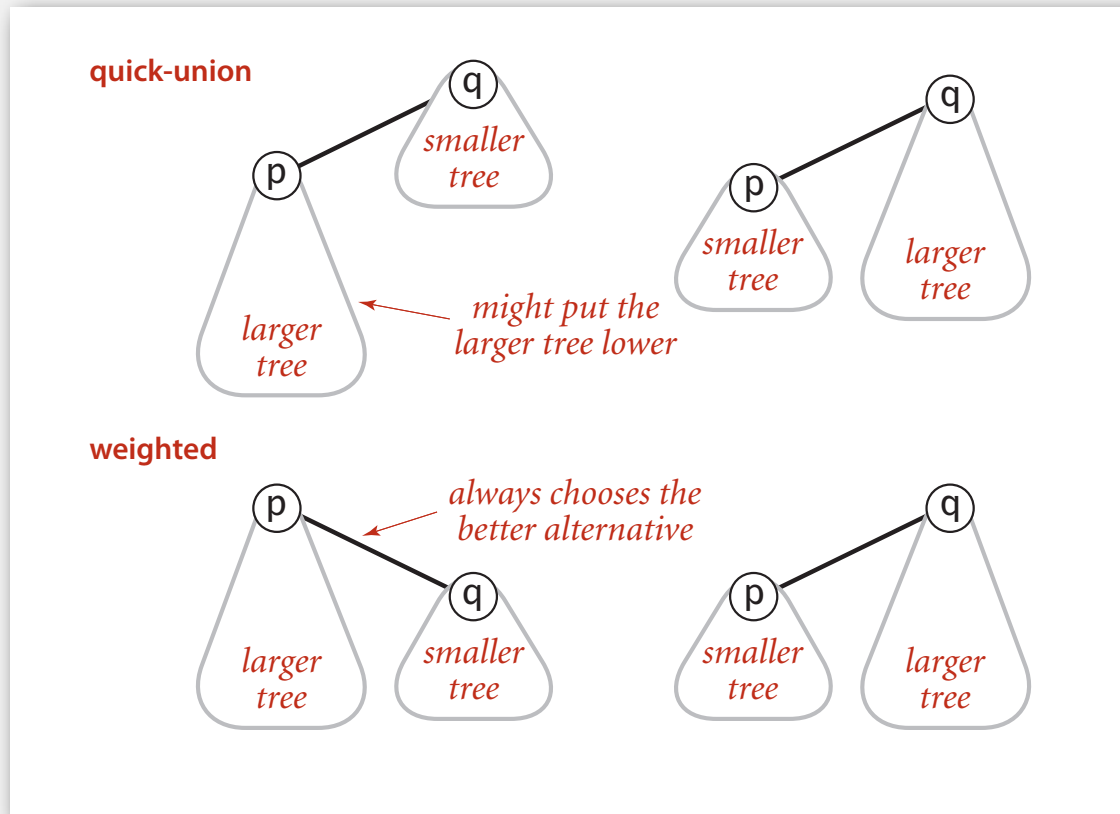
- Trees can get tall.
- Find too expensive (could be N array accesses).

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

Improvement 1: weighting

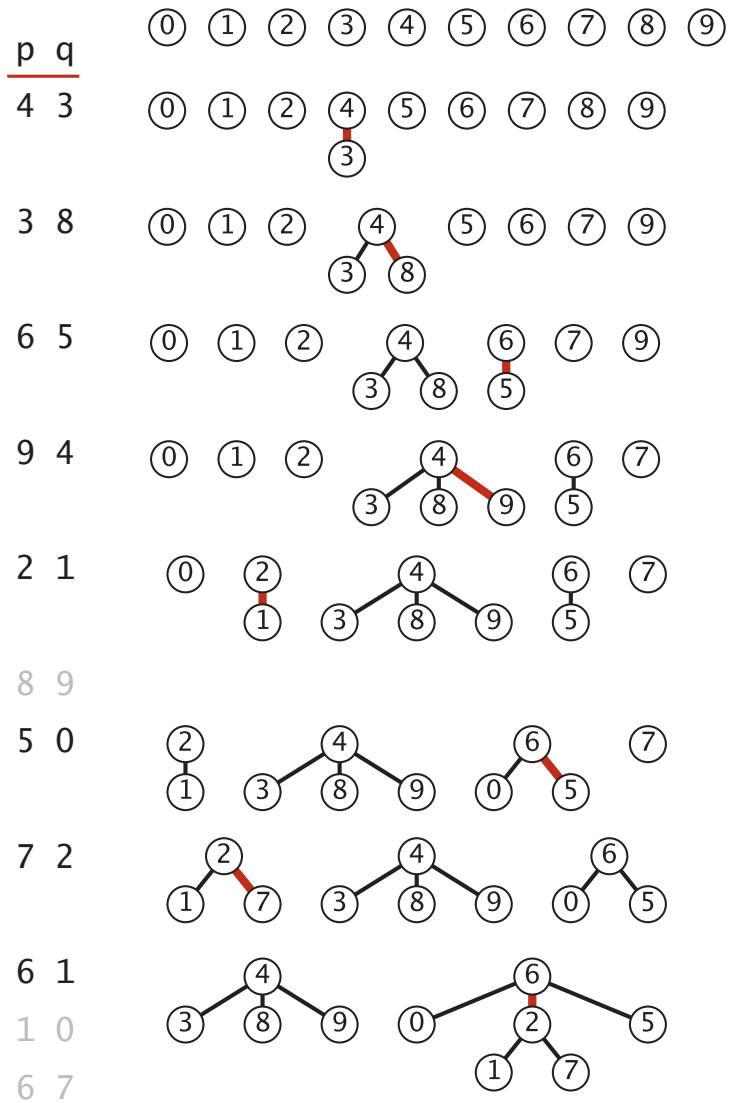
Weighted quick-union.

- Modify quick-union to avoid tall trees.
- Keep track of size of each tree (number of objects).
- Balance by linking small tree below large one.

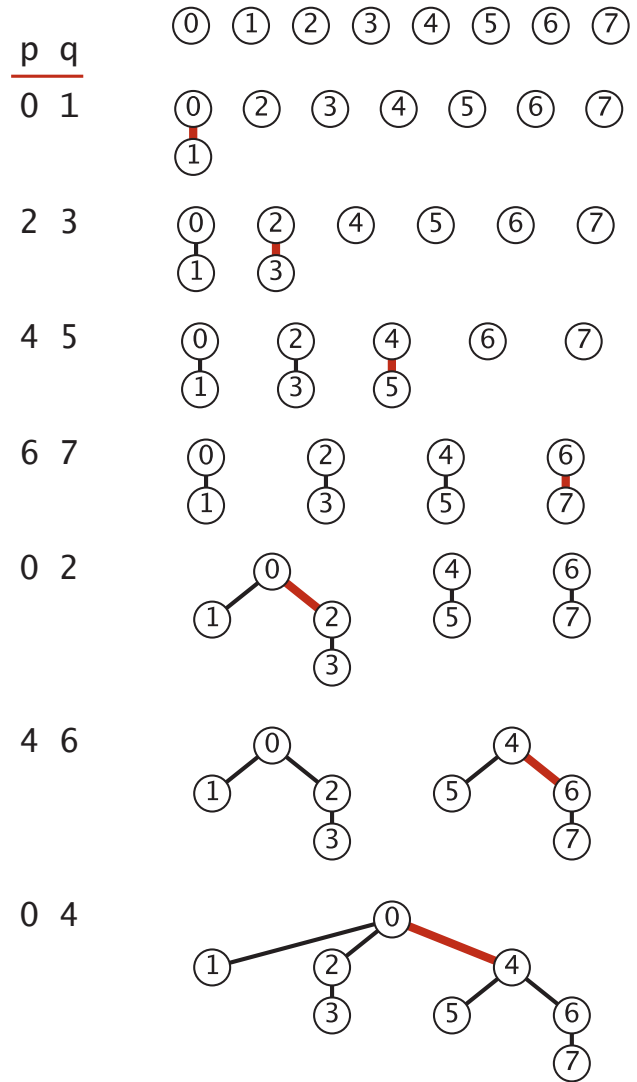


Weighted quick-union examples

reference input

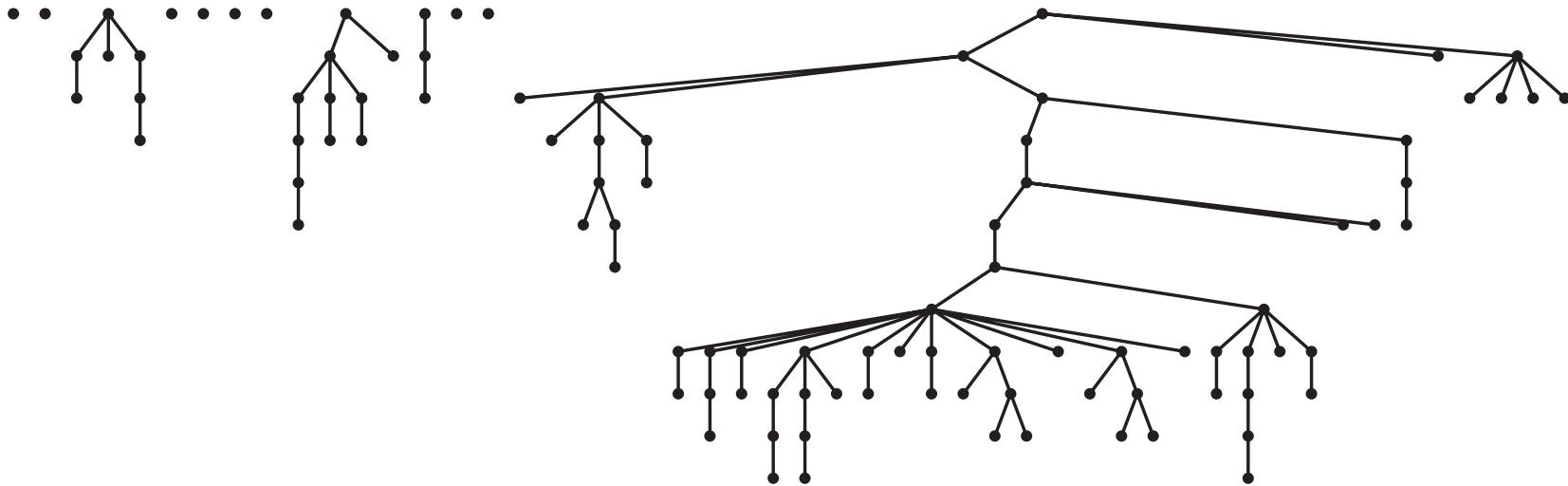


worst-case input



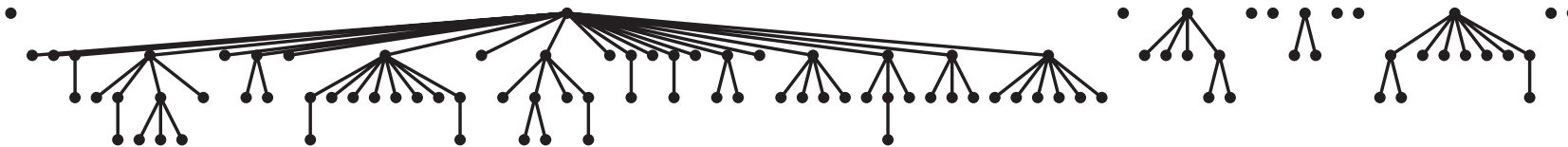
Quick-union and weighted quick-union example

quick-union



average distance to root: 5.11

weighted



average distance to root: 1.52

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. Identical to quick-union.

```
return root(p) == root(q);
```

Union. Modify quick-union to:

- Merge smaller tree into larger tree.
- Update the `sz[]` array.

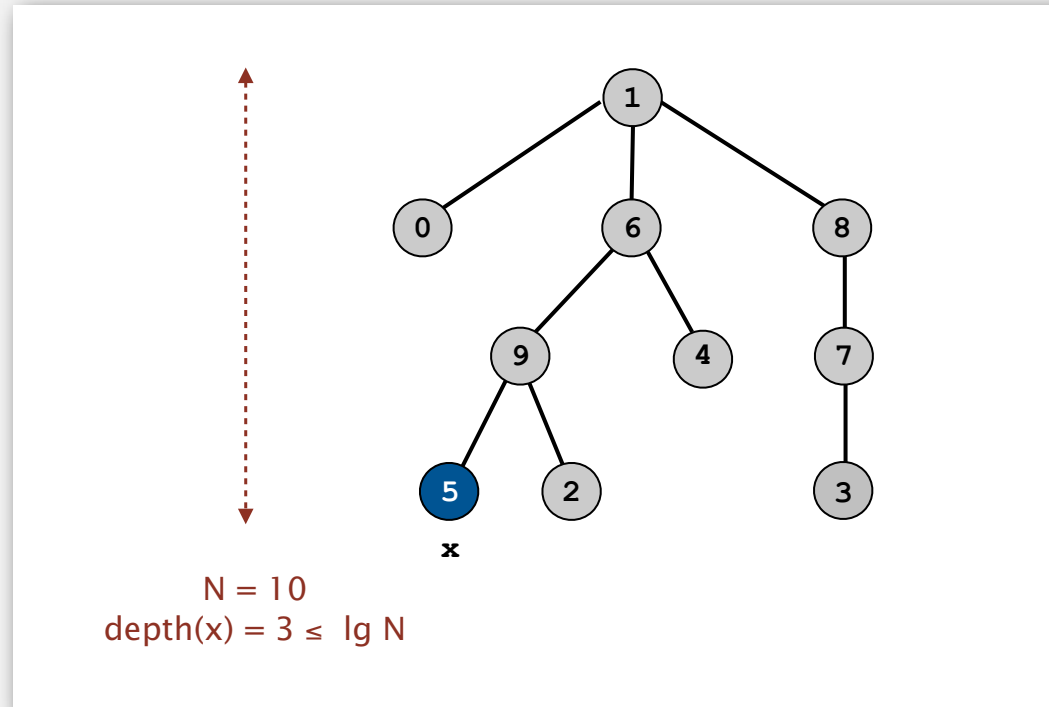
```
int i = root(p);
int j = root(q);
if (sz[i] < sz[j]) { id[i] = j; sz[j] += sz[i]; }
else                { id[j] = i; sz[i] += sz[j]; }
```

Weighted quick-union analysis

Running time.

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

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



Weighted quick-union analysis

Running time.

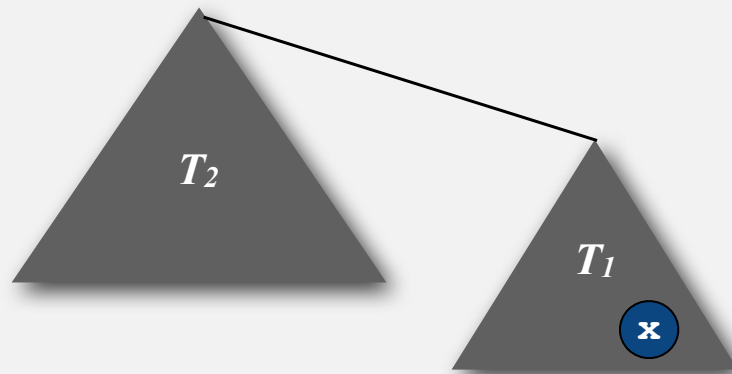
- Find: takes time proportional to depth of p and q .
- Union: takes constant time, given roots.

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

Pf. When does depth of x increase?

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

- The size of the 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 and q .
- Union: takes constant time, given roots.

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

algorithm	init	union	find
quick-find	N	N	1
quick-union	N	N †	N
weighted QU	N	$\lg N$ †	$\lg N$

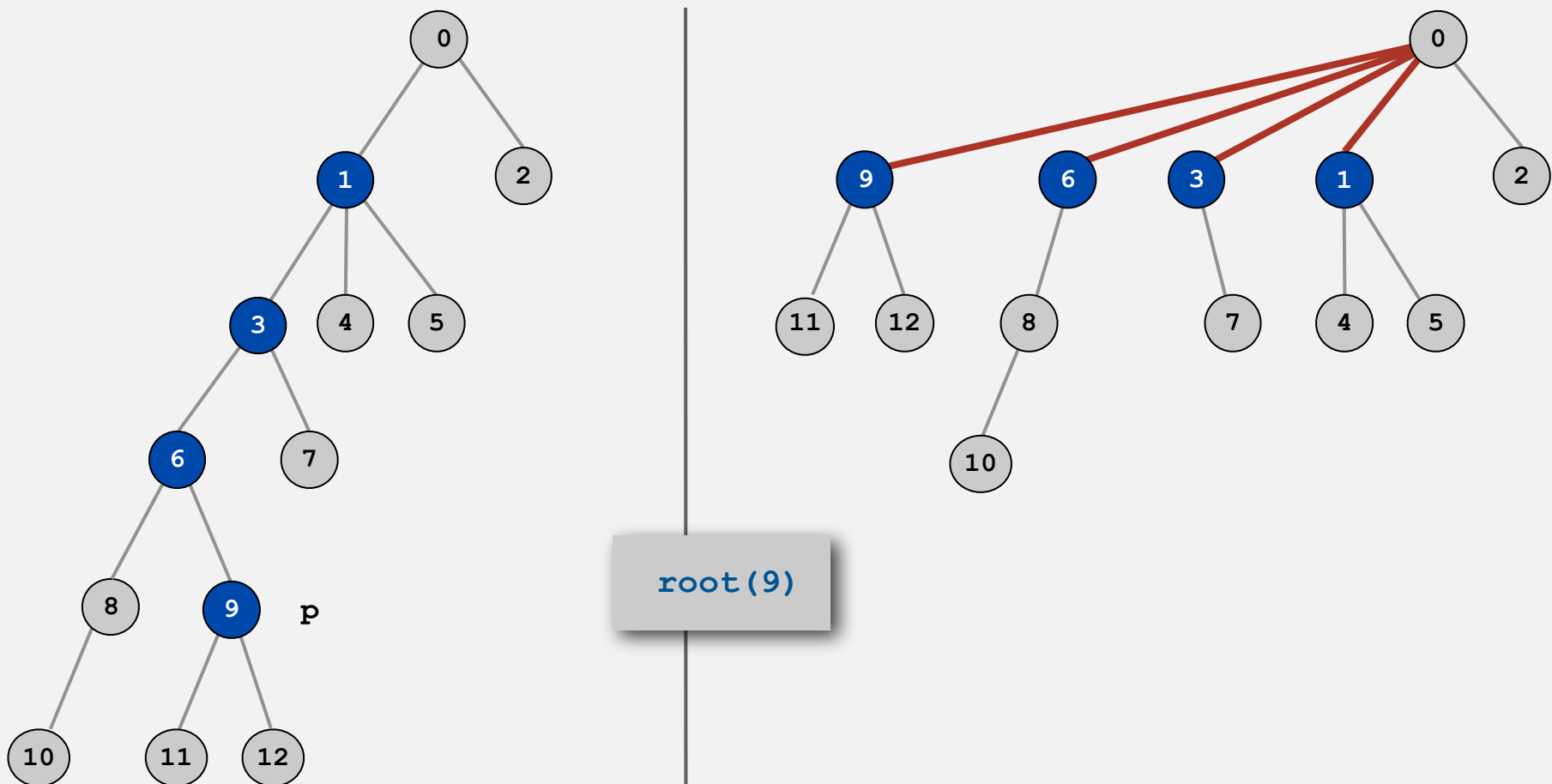
† includes cost of finding root

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.



Path compression: Java implementation

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

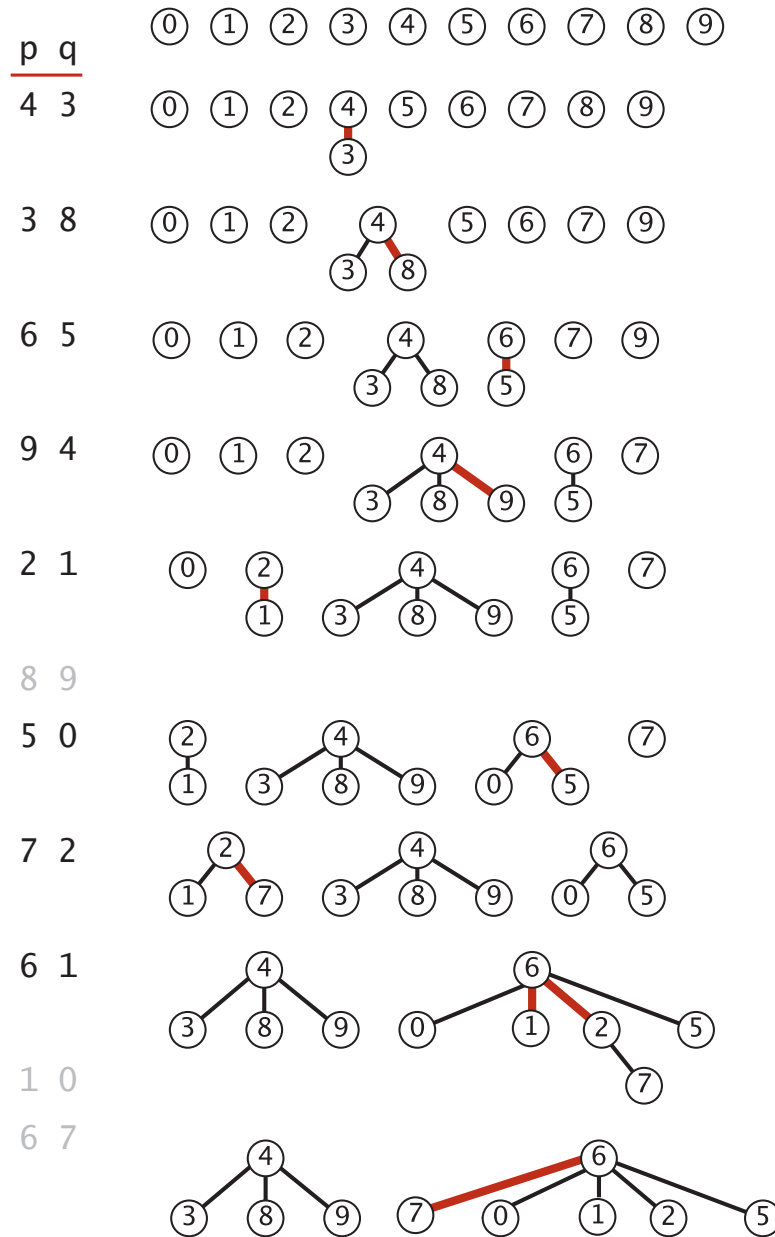
Simpler one-pass variant: halve the path length by making every other node in path point to its grandparent.

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

← only one extra line of code !

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

Weighted quick-union with path compression example



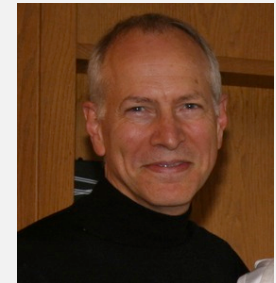
1 linked to 6 because of path compression

7 linked to 6 because of path compression

Weighted quick-union with path compression: amortized analysis

Proposition. Starting from an empty data structure, any sequence of M union-find operations on N objects makes at most proportional to $N + M \lg^* N$ array accesses.

- Proof is very difficult.
- Can be improved to $N + M \alpha(M, N)$. ← see COS 423
- But the algorithm is still simple!



Bob Tarjan
(Turing Award '86)

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.

because $\lg^* N$ is a constant in this universe

Amazing fact. No linear-time algorithm exists.

in "cell-probe" model of computation

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

\lg^* function

Summary

Bottom line. WQUPC 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	$N + M \lg^* N$

M union-find operations on a set of N objects

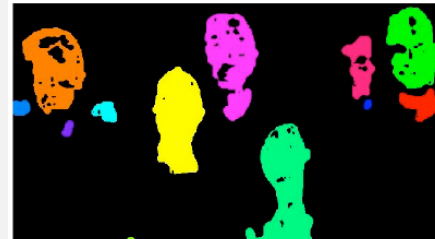
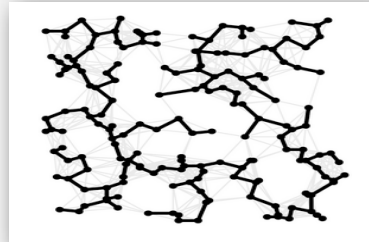
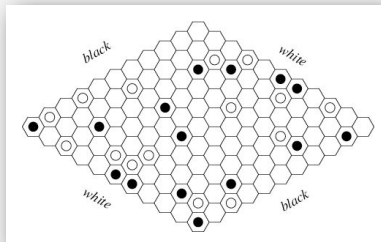
Ex. [10^9 unions and finds with 10^9 objects]

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

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

Union-find applications

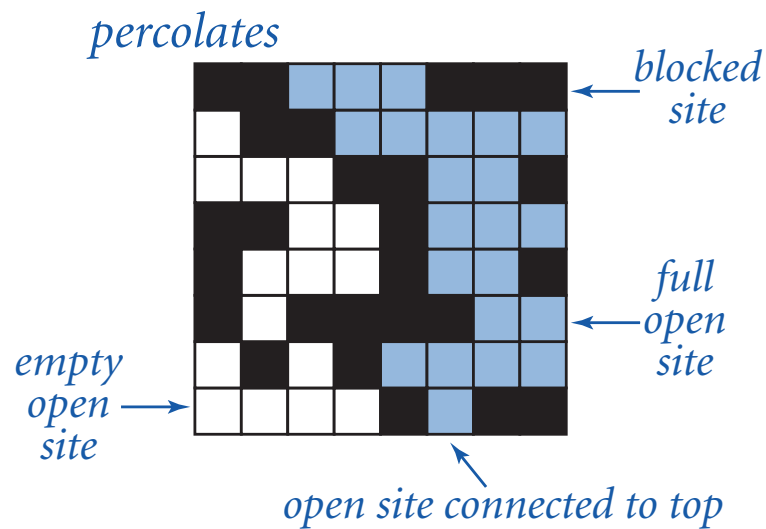
- Percolation.
- Games (Go, Hex).
- ✓ Network connectivity.
- Least common ancestor.
- Equivalence of finite state automata.
- Hoshen-Kopelman algorithm in physics.
- Hinley-Milner polymorphic type inference.
- Kruskal's minimum spanning tree algorithm.
- Compiling equivalence statements in Fortran.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



Percolation

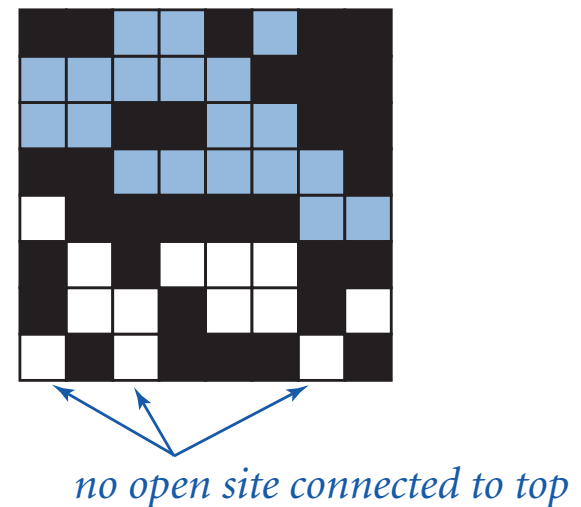
A model for many physical systems:

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



$N = 8$

does not percolate



Percolation

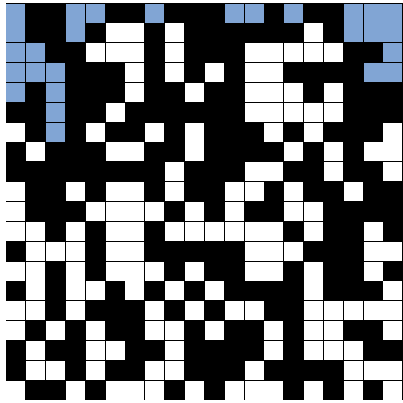
A model for many physical systems:

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

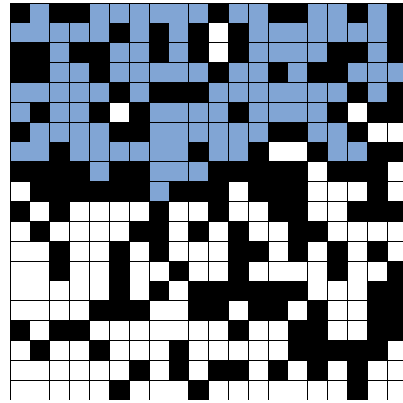
model	system	vacant site	occupied site	percolates
electricity	material	conductor	insulated	conducts
fluid flow	material	empty	blocked	porous
social interaction	population	person	empty	communicates

Likelihood of percolation

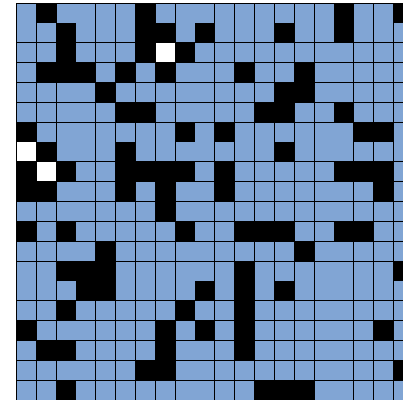
Depends on site vacancy probability p .



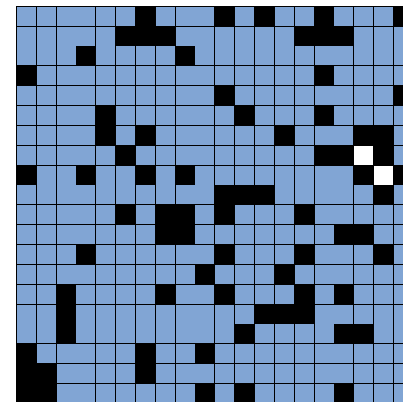
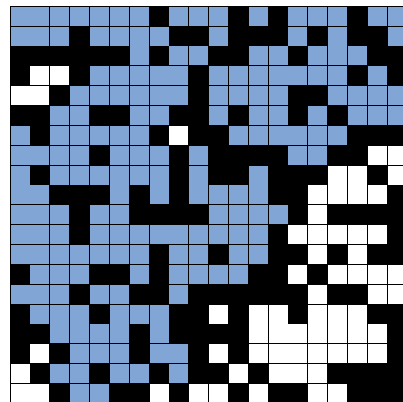
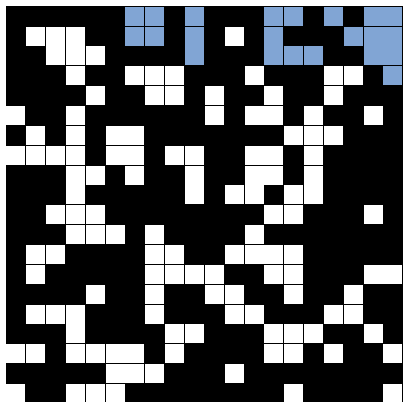
p low (0.4)
does not percolate



p medium (0.6)
percolates?



p high (0.8)
percolates

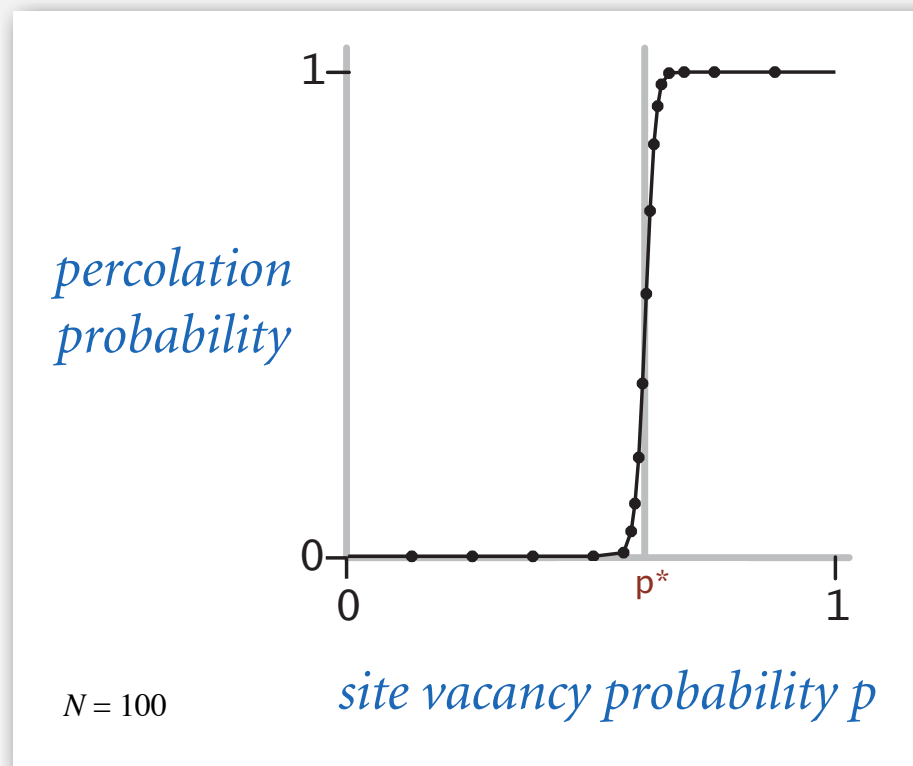


Percolation phase transition

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

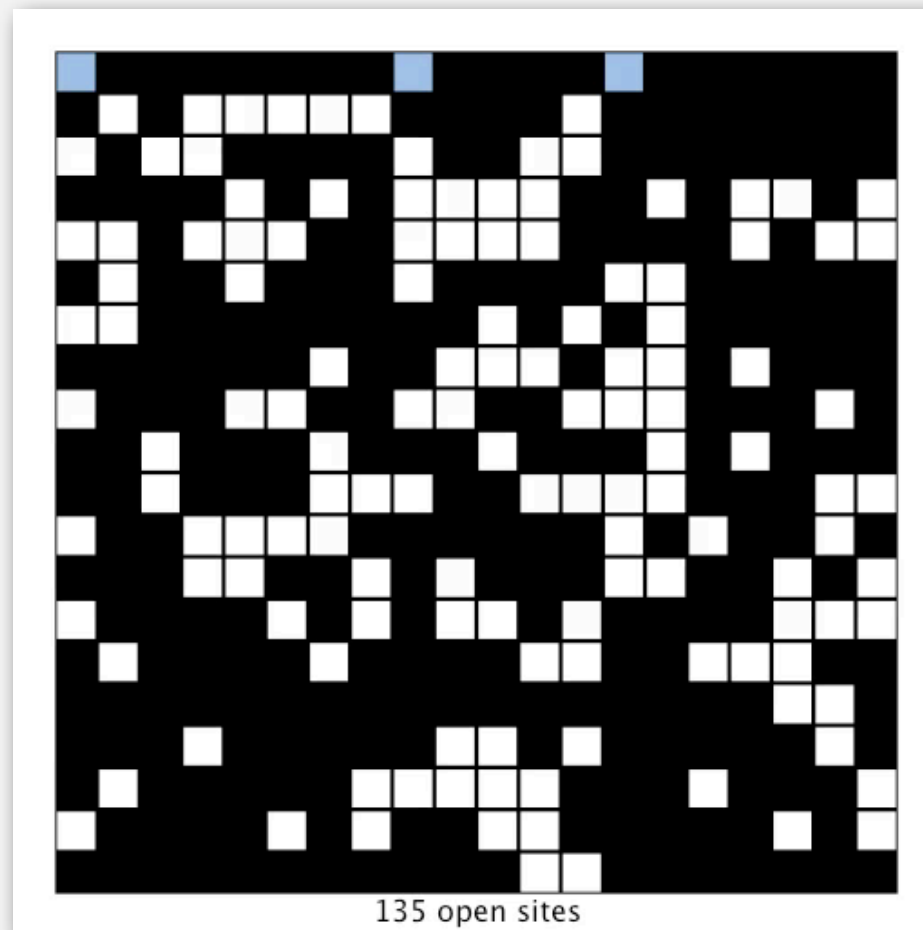
- $p > p^*$: almost certainly percolates.
- $p < p^*$: almost certainly does not percolate.

Q. What is the value of p^* ?

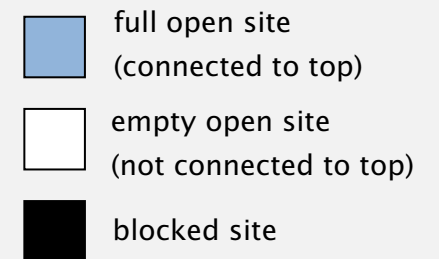


Monte Carlo simulation

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



$N = 20$



UF solution to find percolation threshold

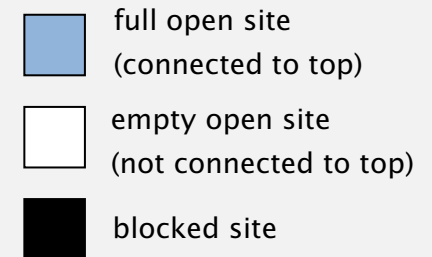
How to check whether system percolates?

- Create an object for each site.
- Sites are in same set if connected by open sites.
- Percolates if any site in top row is in same set as any site in bottom row.

↖ brute force algorithm: check all N^2 pairs

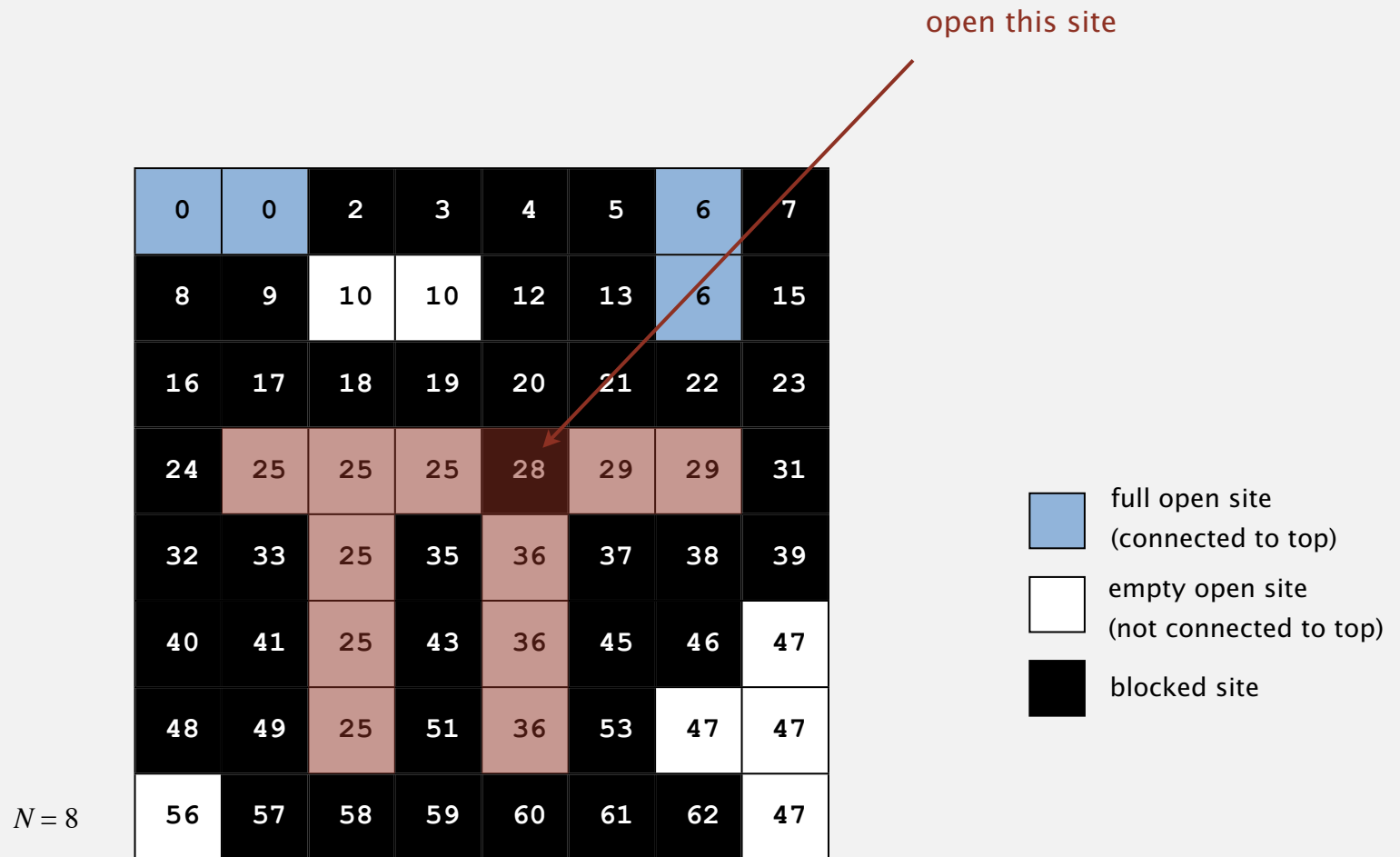
0	0	2	3	4	5	6	7
8	9	10	10	12	13	6	15
16	17	18	19	20	21	22	23
24	25	25	25	28	29	29	31
32	33	25	35	36	37	38	39
40	41	25	43	36	45	46	47
48	49	25	51	36	53	47	47
56	57	58	59	60	61	62	47

$N = 8$



UF solution to find percolation threshold

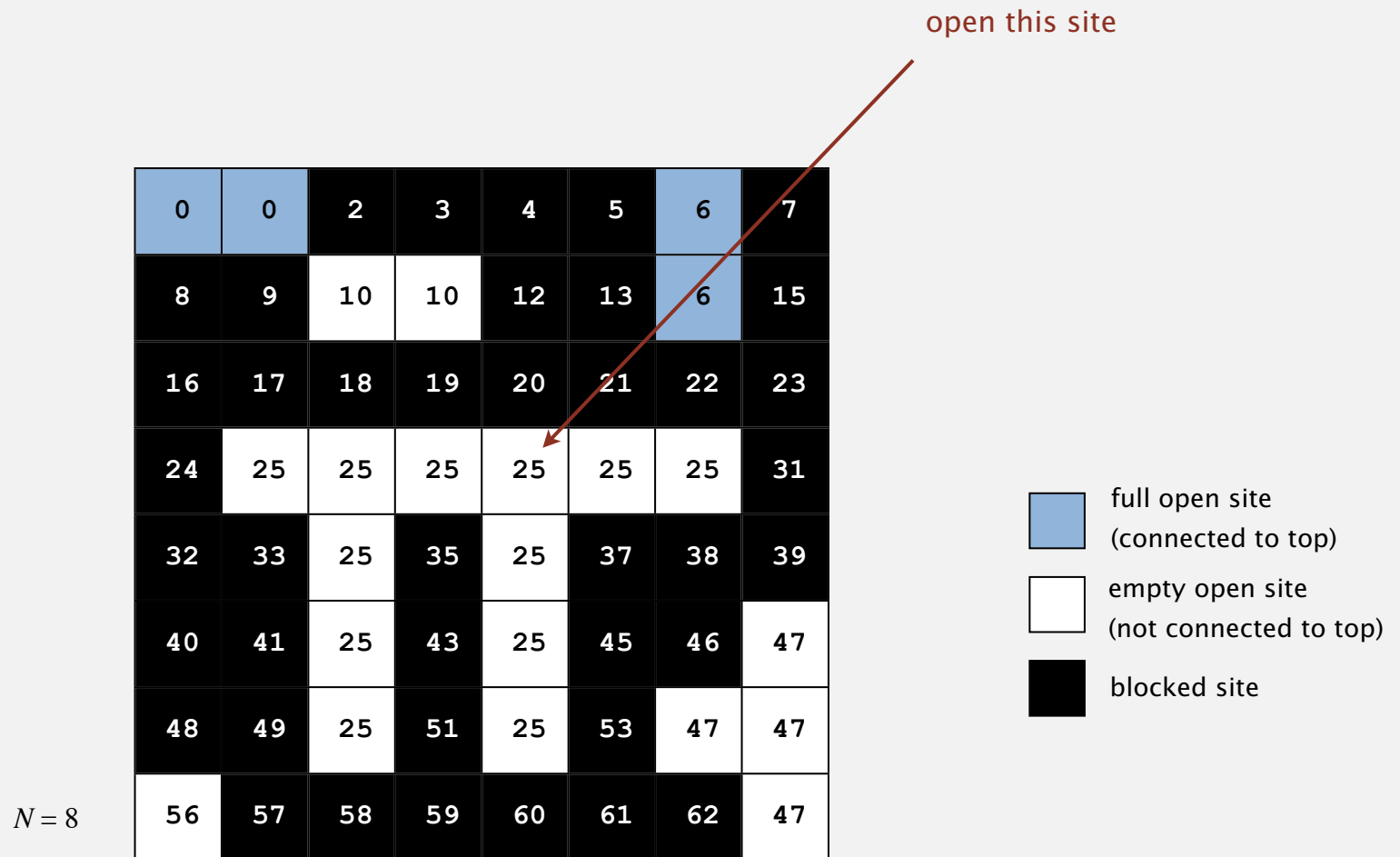
Q. How to declare a new site open?



UF solution to find percolation threshold

Q. How to declare a new site open?

A. Take union of new site and all adjacent open sites.






UF solution: a critical optimization

Q. How to avoid checking all pairs of top and bottom sites?

0	0	2	3	4	5	6	7
8	9	10	10	12	13	6	15
16	17	18	19	20	21	22	23
24	25	25	25	25	25	25	31
32	33	25	35	25	37	38	39
40	41	25	43	25	45	46	47
48	49	25	51	25	53	47	47
56	57	58	59	60	61	62	47

$N = 8$

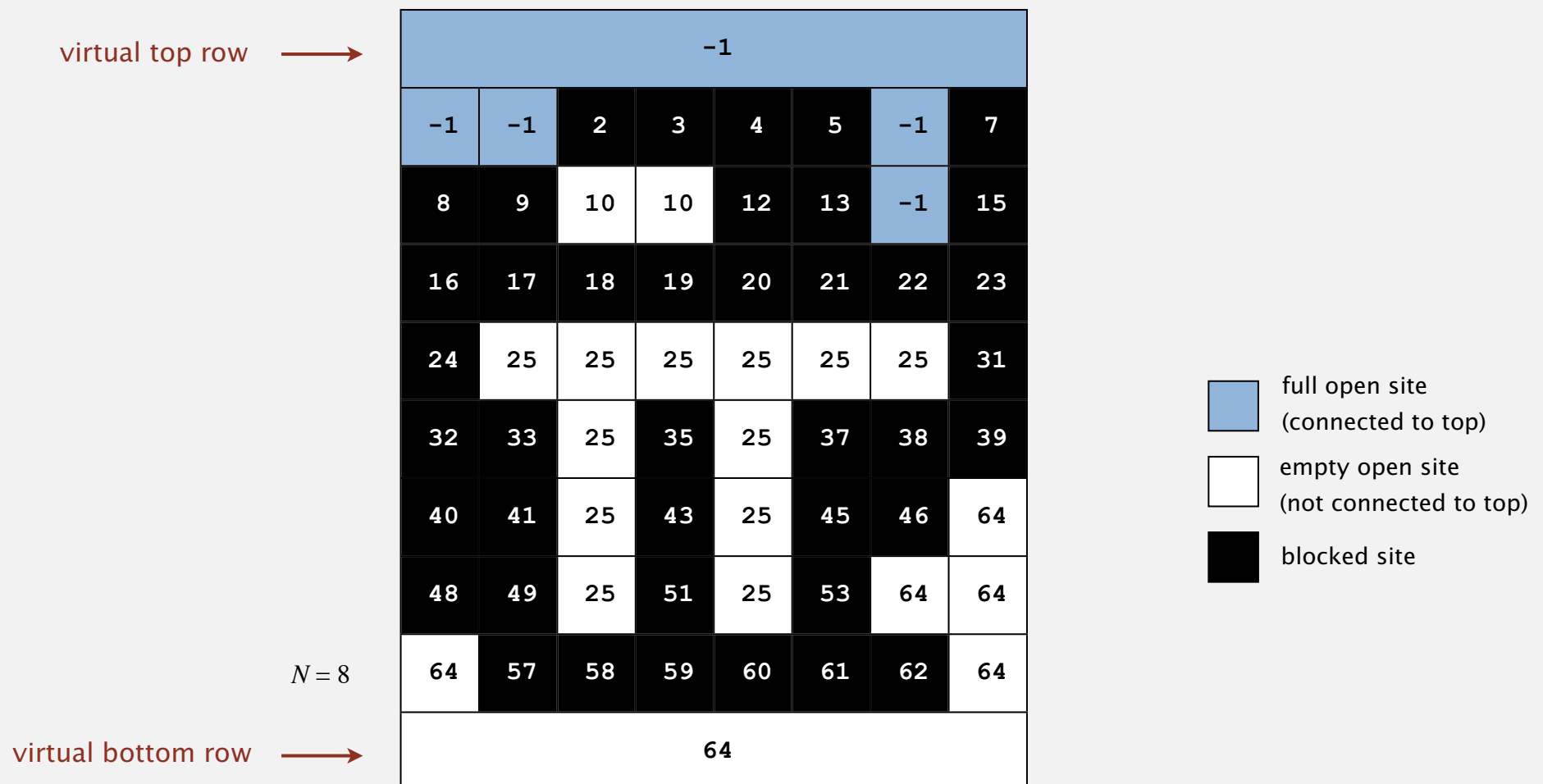
-  full open site
(connected to top)
-  empty open site
(not connected to top)
-  blocked site

UF solution: a critical optimization

Q. How to avoid checking all pairs of top and bottom sites?

A. Create a virtual top and bottom objects;

system percolates when virtual top and bottom objects are in same set.

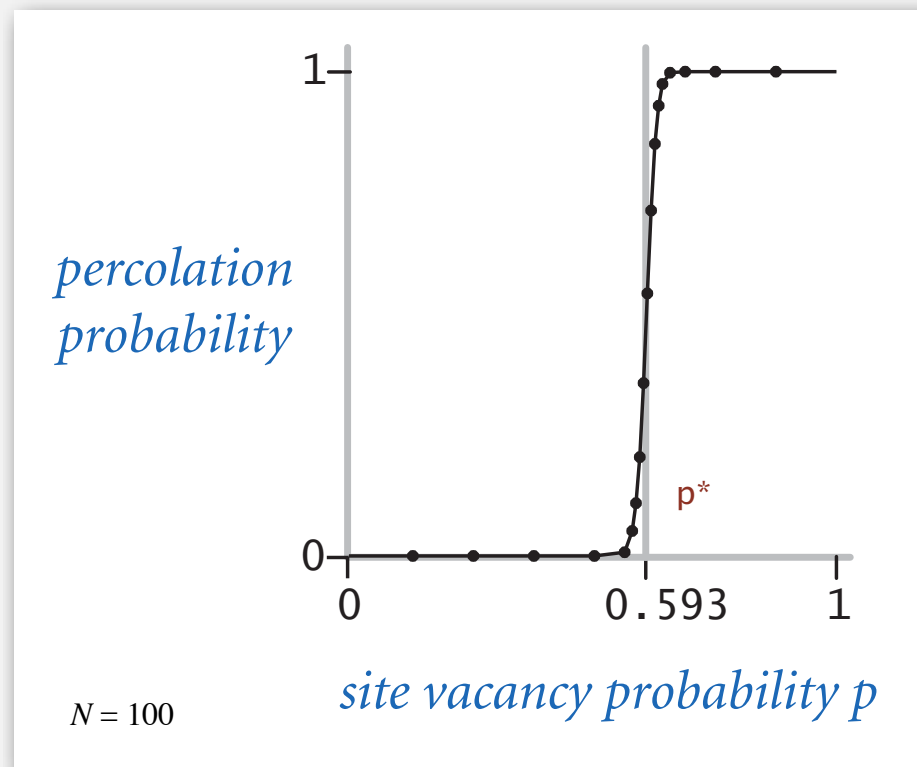


Percolation threshold

Q. What is percolation threshold p^* ?

A. About 0.592746 for large square lattices.

↑
percolation constant known
only via simulation



Fast algorithm **enables** accurate answer to scientific question.

Subtext 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.