# COS 226

## Algorithms and Data Structures
## Princeton University
## Spring 2010

## Robert Sedgewick

# 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, ...

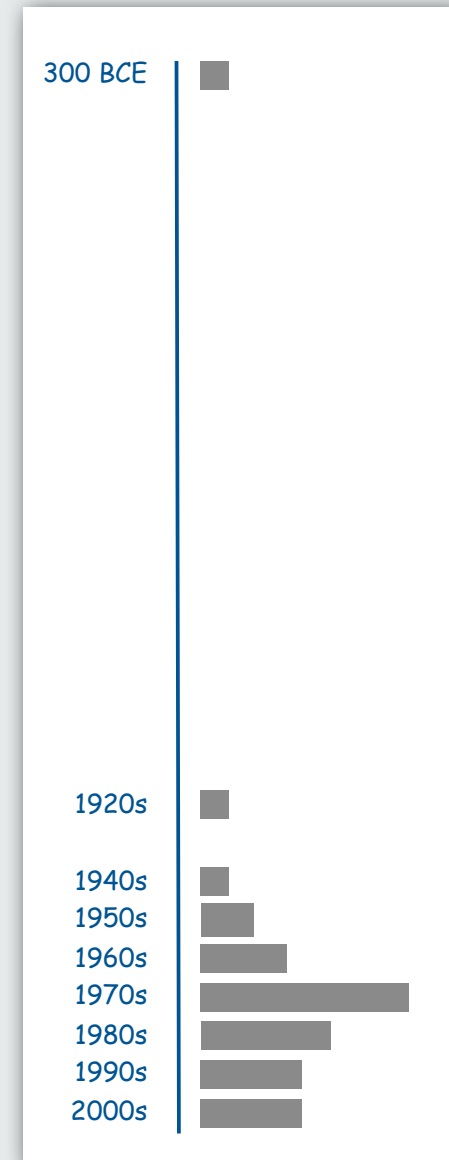Transportation.  Airline crew scheduling, map routing, ...

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

...

## Why study algorithms?
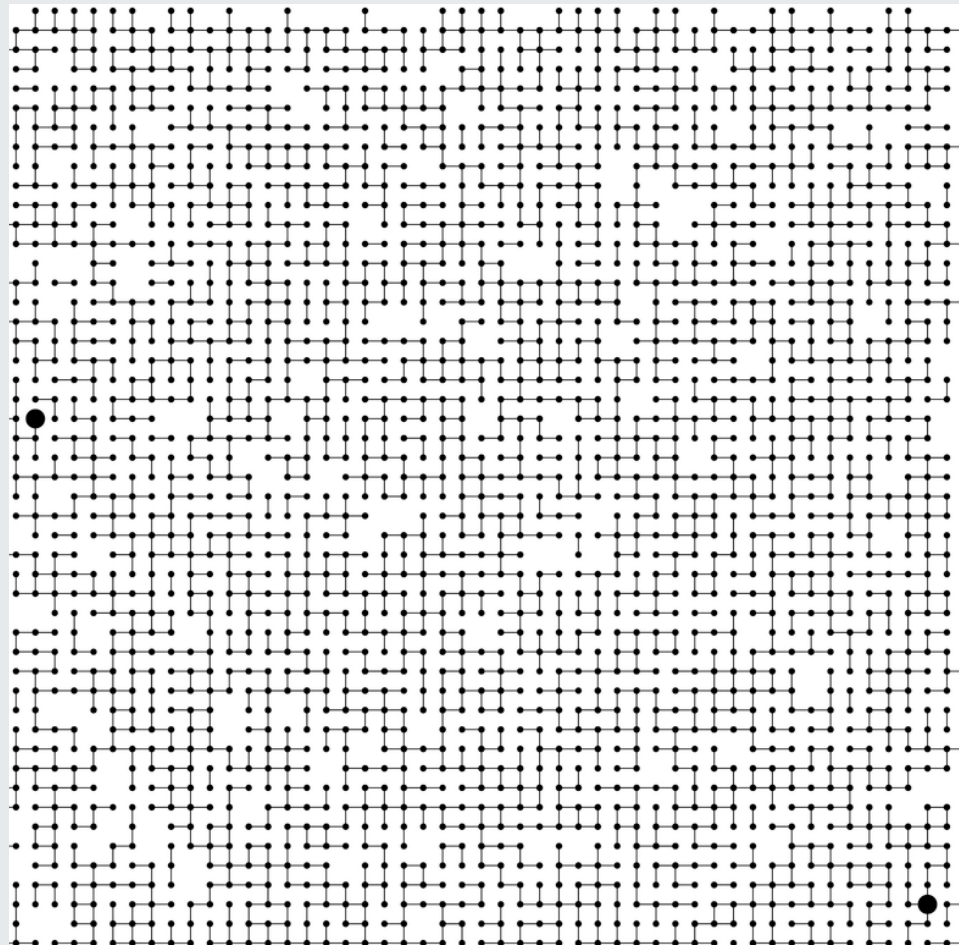
**Old roots, new opportunities.**

- 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.* " — Francis Sullivan

" *An algorithm must be seen to be believed.* " — D. E. Knuth

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

Computational models are replacing mathematical models in scientific inquiry.

$$E = mc^2$$

$$F = ma$$

$$F = \frac{Gm_1m_2}{r^2}$$

$$\left[-\frac{h^2}{2m}\nabla^2 + V(r)\right]\Psi(r) = E\,\Psi(r)$$

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."* — Avi Wigderson

# Why study algorithms?

For fun and profit.

## 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?

## Coursework and grading

8 programming assignments.  45%

- Electronic submission.
- Due 11pm, starting Wednesay 9/23.
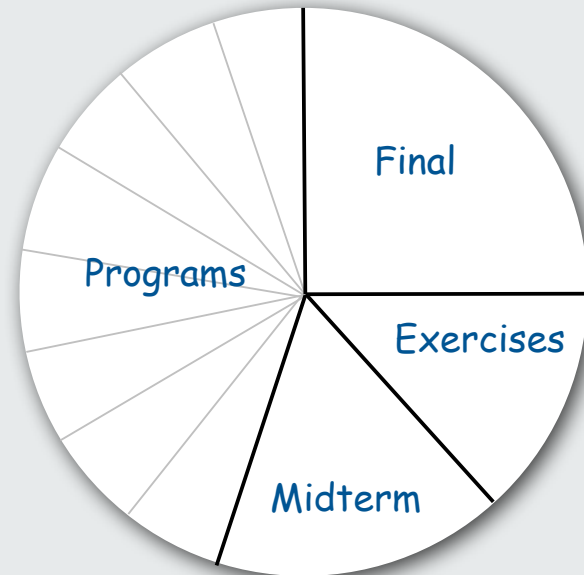
Exercises.  15%

- Due in lecture, starting Tuesday 9/22.

Exams.

- Closed-book with cheatsheet.
- Midterm.  15%
- Final.       25%

Staff discretion.  To adjust borderline cases.

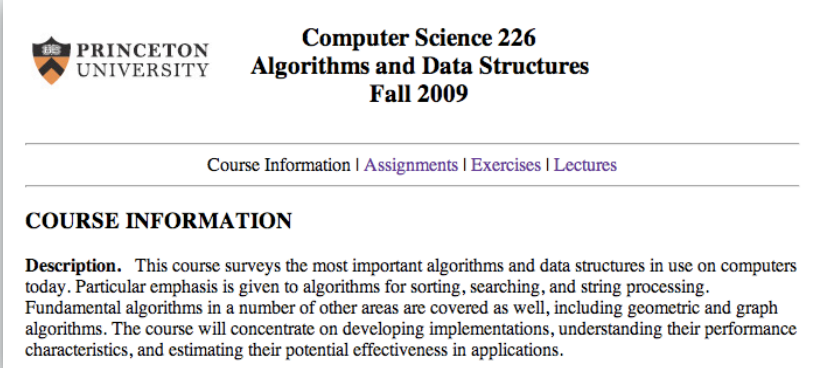everyone needs to meet me in office hours

## Resources (web)

### Course content.

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



`http://www.princeton.edu/~cos226`

### Booksites.

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



`http://www.cs.princeton.edu/IntroProgramming`
`http://www.cs.princeton.edu/algs4`

# 1.5 Case Study



‣ 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

3

# 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
```

# Network connectivity: larger example

**Q.** Is there a path from p to q?



p

q

**A.** Yes.  ← but finding the path is more difficult: stay tuned (Chapter 4)

## Modeling the objects

Dynamic connectivity applications involve manipulating objects of all types.

- Variable name aliases.
- Pixels in a digital photo.
- Computers in a network.
- Web pages on the Internet.
- Transistors in a computer chip.
- 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)

## Modeling the connections

Transitivity. 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.



{ 1  5  6 }  { 2  3  4  7 }   { 0  8 }

connected components

# Implementing the operations

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

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



union(4, 8)

{ 1 5 6 } { 2 3 4 7 }  { 0 8 }          { 1 5 6 } { 0 2 3 4 7 8 }

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.
- Find queries and union commands may be intermixed.

```
public class UnionFind
```

| | |
|---|---|
| `UnionFind(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 set?* |
| `void unite(int p, int q)` | *replace sets containing p and q with their union* |

# Quick-find [eager approach]

## Data structure.

- Integer array `id[]` of size `N`.
- Interpretation: `p` and `q` are connected if 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

# Quick-find  [eager approach]

Data structure.

- Integer array `id[]` of size `N`.
- Interpretation:  `p` and `q` are connected if 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 not connected

## Quick-find  [eager approach]

Data structure.

- Integer array `id[]` of size `N`.
- Interpretation:  `p` and `q` are connected if 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 not connected

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 |

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

problem: many values can change

# Quick-find example

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 3–4 | 0 | 1 | 2 | 4 | 4 | 5 | 6 | 7 | 8 | 9 |
| 4–9 | 0 | 1 | 2 | 9 | 9 | 5 | 6 | 7 | 8 | 9 |
| 8–0 | 0 | 1 | 2 | 9 | 9 | 5 | 6 | 7 | 0 | 9 |
| 2–3 | 0 | 1 | 9 | 9 | 9 | 5 | 6 | 7 | 0 | 9 |
| 5–6 | 0 | 1 | 9 | 9 | 9 | 6 | 6 | 7 | 0 | 9 |
| 5–9 | 0 | 1 | 9 | 9 | 9 | 9 | 9 | 7 | 0 | 9 |
| 7–3 | 0 | 1 | 9 | 9 | 9 | 9 | 9 | 9 | 0 | 9 |
| 4–8 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6–1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |

problem: many values can change

## Quick-find: Java implementation

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

    public QuickFind(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 unite(int p, int q)
    {
        int pid = id[p];
        for (int i = 0; i < id.length; i++)
            if (id[i] == pid) id[i] = id[q];
    }
}
```

set id of each object to itself
(N operations)

check if `p` and `q` have same id
(1 operation)

change all entries with `id[p]` to `id[q]`
(N operations)

## Quick-find is too slow

Quick-find defect.

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

| algorithm | union | find |
|-----------|-------|------|
| quick-find | N | 1 |

Ex.  Takes $N^2$ operations 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*

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

```
   i   0  1  2  3  4  5  6  7  8  9
id[i]  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 not connected

# 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*

```
  i    0  1  2  3  4  5  6  7  8  9
id[i]  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 not connected

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

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

*only one value changes*

21

# Quick-union example



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

problem:
trees can get tall

# Quick-union:  Java implementation

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

   public QuickUnion(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 unite(int p, int q)
   {
      int i = root(p), j = root(q);
      id[i] = j;
   }
}
```

set id of each object to itself
(N operations)

chase parent pointers until reach root
(depth of i operations)

check if p and q have same root
(depth of p and q operations)

change root of p to point to root of q
(depth of p and q operations)

# Quick-union is also too slow

## Quick-find defect.

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

## Quick-union defect.

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

| algorithm | union | find |
|-----------|-------|------|
| quick-find | N | 1 |
| quick-union | N † | N |

← worst case

† includes cost of finding root

‣ 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 set.
- Balance by linking small tree below large one.

Ex.  Union of 3 and 5.

- Quick union:  link 9 to 6.
- Weighted quick union:  link 6 to 9.

# Weighted quick-union example



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

no problem:
trees stay flat

27

# 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

Analysis.

- 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.



N = 10
depth(x) = 3 ≤ lg N

## Weighted quick-union analysis

Analysis.

- 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

Analysis.

- 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 | union | find |
|-----------|-------|------|
| quick-find | N | 1 |
| quick-union | N † | N |
| weighted QU | 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 `root(p)`.



root(9)

Path compression:  Java implementation

Standard implementation:  add second loop to `root()` 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.

```java
public int root(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 example



3–4    0  1  2  3  3  5  6  7  8  9

4–9    0  1  2  3  3  5  6  7  8  3

8–0    8  1  2  3  3  5  6  7  8  3

2–3    8  1  3  3  3  5  6  7  8  3

5–6    8  1  3  3  3  5  5  7  8  3

5–9    8  1  3  3  3  3  5  7  8  3

7–3    8  1  3  3  3  3  5  3  8  3

4–8    8  1  3  3  3  3  5  3  3  3

6–1    8  3  3  3  3  3  3  3  3  3

no problem:
trees stay VERY flat

## WQUPC performance

Proposition. [Tarjan 1975] Starting from an empty data structure, any sequence of M union and find ops on N objects takes $O(N + M \lg^* N)$ time.

- Proof is very difficult.
- But the algorithm is still simple!

actually $O(N + M \, \alpha(M, N))$
see COS 423

Linear algorithm?

- 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

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

lg* function
number of times needed to take
the lg of a number until reaching 1

Amazing fact. No linear-time linking strategy exists.

## 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.



*percolates*       *does not percolate*

*blocked site*

*full open site*

*empty open site*

*site connected to top*

*no open site connected to top*

*N = 8*

## 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*
*does not percolate*

*p medium*
*percolates?*

*p high*
*percolates*

$N = 20$

# 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* ?



percolation probability

$N = 100$

site vacancy probability 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*.



Sites = 135

full open site
(connected to top)

empty open site
(not connected to top)

blocked site

# 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 needs to check $N^2$ 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 |

$N = 8$

■ full open site (connected to top)

□ empty open site (not connected to top)

■ blocked site

44

# UF solution to find percolation threshold

Q. How to declare a new site open?

open this site

| 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 |

$N = 8$

full open site
(connected to top)

empty open site
(not connected to top)

blocked site

# 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.

open this site

| 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

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

$N = 8$

# 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.



virtual top row →

| 0 | 0 | 2 | 3 | 4 | 5 | 0 | 7 |
| 8 | 9 | 10 | 10 | 12 | 13 | 0 | 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 |
| 47 | 57 | 58 | 59 | 60 | 61 | 62 | 47 |

virtual top row: 0

virtual bottom row → 47

full open site (connected to top)

empty open site (not connected to top)

blocked site

$N = 8$

48

## Percolation threshold

Q. What is percolation threshold p* ?

A. About 0.592746 for large square lattices.

↑
percolation constant known
only via simulation

## 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.

# 1.4  Analysis of Algorithms



▸ estimating running time

▸ mathematical analysis

▸ order-of-growth hypotheses

▸ input models

▸ measuring space

*Reference:  Intro to Programming in Java, Section 4.1*

# Cast of characters

**Programmer** needs to develop a working solution.

**Client** wants problem solved efficiently.

**Theoretician** wants to understand.

**Student** might play any or all of these roles someday.

Basic **blocking and tackling** is sometimes necessary. [this lecture]

" *As soon as an Analytic Engine exists, it will necessarily guide the future course of the science. Whenever any result is sought by its aid, the question will arise—By what course of calculation can these results be arrived at by the machine in the shortest time?* " — Charles Babbage



Charles Babbage (1864)



Analytic Engine

how many times do you have to turn the crank?

# Reasons to analyze algorithms

Predict performance.

Compare algorithms.

this course (COS 226)

Provide guarantees.

Understand theoretical basis.          theory of algorithms (COS 423)

Primary practical reason: avoid performance bugs.



client gets poor performance because programmer
did not understand performance characteristics

## Some algorithmic successes

### Discrete Fourier transform.

- Break down waveform of N samples into periodic components.
- Applications:  DVD, JPEG, MRI, astrophysics, ….
- Brute force:  $N^2$ steps.
- FFT algorithm:  N log N steps, enables new technology.

Friedrich Gauss
1805

## Some algorithmic successes

### N-body Simulation.

- Simulate gravitational interactions among N bodies.
- Brute force:  $N^2$ steps.
- Barnes-Hut:  N log N steps, enables new research.

Andrew Appel
PU '81



time

64T — quadratic

32T —

16T — linearithmic

8T — linear

size → 1K 2K 4K 8K



Galaxies NGC 2207 and IC 2163

Hubble Heritage

‣ **estimating running time**

‣ mathematical analysis

‣ order-of-growth hypotheses

‣ input models

‣ measuring space

## Scientific analysis of algorithms

A framework for predicting performance and comparing algorithms.

### Scientific method.

- Observe some feature of the universe.
- Hypothesize a model that is consistent with observation.
- Predict events using the hypothesis.
- Verify the predictions by making further observations.
- Validate by repeating until the hypothesis and observations agree.

### Principles.

- Experiments must be reproducible.
- Hypotheses must be falsifiable.

### Universe = computer itself.

## Experimental algorithmics

Every time you run a program you are doing an experiment!

*Why is my program so slow ??*

First step.  Debug your program!

Second step.  Choose input model for experiments.

Third step.  Run and time the program for problems of increasing size.

## Example: 3-sum

3-sum. Given $N$ integers, find all triples that sum to exactly zero.

```
% more input8.txt
8
 30 -30 -20 -10 40 0 10 5

% java ThreeSum < input8.txt
  4
  30 -30   0
  30 -20 -10
-30 -10  40
-10   0  10
```

Context. Deeply related to problems in computational geometry.

## 3-sum: brute-force algorithm

```java
public class ThreeSum
{
   public static int count(int[] a)
   {
      int N = a.length;
      int cnt = 0;
      for (int i = 0; i < N; i++)
         for (int j = i+1; j < N; j++)
            for (int k = j+1; k < N; k++)          ← check each triple
               if (a[i] + a[j] + a[k] == 0)        ← ignore overflow
                  cnt++;
      return cnt;
   }

   public static void main(String[] args)
   {
      long[] a = StdArrayIO.readInt1D();
      StdOut.println(count(a));
   }
}
```

## Empirical analysis

Run the program for various input sizes and measure running time.

**ThreeSum.java**

| N | time (seconds) [†] |
|---|---|
| 1000 | 0.26 |
| 2000 | 2.16 |
| 4000 | 17.18 |
| 8000 | 137.76 |

† Running Linux on Sun-Fire-X4100

# Measuring the running time

Q. How to time a program?

A. Manual.

```
% java ThreeSum < 1Kints.txt
```

tick tick tick

0

```
% java ThreeSum < 2Kints.txt
```

tick tick tick tick tick tick
tick tick tick tick tick tick
tick tick tick tick tick tick
tick tick tick tick tick tick

2
391930676 -763182495 371251819
-326747290 802431422 -475684132

## Measuring the running time

Q. How to time a program?

A. Automatic.

```
Stopwatch stopwatch = new Stopwatch();

ThreeSum.count(a);

double time = stopwatch.elapsedTime();
StdOut.println("Running time: " + time + " seconds");
```

client code

```
public class Stopwatch
{
    private final long start = System.currentTimeMillis();

    public double elapsedTime()
    {
        long now = System.currentTimeMillis();
        return (now - start) / 1000.0;
    }
}
```

implementation (part of `stdlib.jar`, see `http://www.cs.princeton.edu/introcs/stdlib`)

## Data analysis

Plot running time as a function of input size N.

# Data analysis

Log-log plot. Plot running time vs. input size $N$ on log-log scale.



Regression. Fit straight line through data points: $a N^b$. ← power law

Hypothesis. Running time grows with the cube of the input size: $a N^3$. ← slope

## Doubling hypothesis

**Doubling hypothesis.** Quick way to estimate $b$ in a power law hypothesis.

Run program, doubling the size of the input.

| N | time (seconds) [†] | ratio | lg ratio |
|---|---|---|---|
| 500 | 0.03 | - | |
| 1,000 | 0.26 | 7.88 | 2.98 |
| 2,000 | 2.16 | 8.43 | 3.08 |
| 4,000 | 17.18 | 7.96 | 2.99 |
| 8,000 | 137.76 | 7.96 | 2.99 |

seems to converge to a constant b ≈ 3

**Hypothesis.** Running time is about $a N^b$ with $b = $ lg ratio.

**Caveat.** Can't identify logarithmic factors with doubling hypothesis.

## Prediction and verification

**Hypothesis.** Running time is about $a N^3$ for input of size $N$.

**Q.** How to estimate $a$?

**A.** Run the program!

| N | time (seconds) |
|---|---|
| 4,000 | 17.18 |
| 4,000 | 17.15 |
| 4,000 | 17.17 |

$17.17 = a \times 4000^3$
$\Rightarrow a = 2.7 \times 10^{-10}$

**Refined hypothesis.** Running time is about $2.7 \times 10^{-10} \times N^3$ seconds.

**Prediction.** 1,100 seconds for $N = 16,000$.

**Observation.**

| N | time (seconds) |
|---|---|
| 16384 | 1118.86 |

validates hypothesis!

## Experimental algorithmics

Many obvious factors affect running time:

- Machine.
- Compiler.
- Algorithm.
- Input data.

More factors (not so obvious):

- Caching.
- Garbage collection.
- Just-in-time compilation.
- CPU use by other applications.

**Bad news.** It is often difficult to get precise measurements.

**Good news.** Easier than other sciences.

e.g., can run huge number of experiments

Q. How long does this program take as a function of N?

```
public class EditDistance
{
    String s = StdIn.readString();
    int N = s.length();
    ...
      for (int i = 0; i < N; i++)
         for (int j = 0; j < N; j++)
             distance[i][j] = ...

    ...
}
```

Jenny. ~ $c_1 N^2$ seconds.

Kenny. ~ $c_2 N$ seconds.

| N | time |
|---|------|
| 1,000 | 0.11 |
| 2,000 | 0.35 |
| 4,000 | 1.6 |
| 8,000 | 6.5 |

Jenny

| N | time |
|---|------|
| 250 | 0.5 |
| 500 | 1.1 |
| 1,000 | 1.9 |
| 2,000 | 3.9 |

Kenny

‣ estimating running time
‣ **mathematical analysis**
‣ order-of-growth hypotheses
‣ input models
‣ measuring space

## Mathematical models for running time

**Total running time:** sum of cost × frequency for all operations.

- Need to analyze program to determine set of operations.
- Cost depends on machine, compiler.
- Frequency depends on algorithm, input data.



Donald Knuth
1974 Turing Award

**In principle,** accurate mathematical models are available.

# Cost of basic operations

| operation | example | nanoseconds [†] |
|---|---|---|
| integer add | `a + b` | 2.1 |
| integer multiply | `a * b` | 2.4 |
| integer divide | `a / b` | 5.4 |
| floating point add | `a + b` | 4.6 |
| floating point multiply | `a * b` | 4.2 |
| floating point divide | `a / b` | 13.5 |
| sine | `Math.sin(theta)` | 91.3 |
| arctangent | `Math.atan2(y, x)` | 129.0 |
| … | … | … |

† Running OS X on Macbook Pro 2.2GHz with 2GB RAM

## Cost of basic operations

| operation | example | nanoseconds [†] |
|---|---|---|
| variable declaration | `int a` | $c_1$ |
| assignment statement | `a = b` | $c_2$ |
| integer compare | `a < b` | $c_3$ |
| array element access | `a[i]` | $c_4$ |
| array length | `a.length` | $c_5$ |
| 1D array allocation | `new int[N]` | $c_6\ N$ |
| 2D array allocation | `new int[N][N]` | $c_7\ N^2$ |
| string length | `s.length()` | $c_8$ |
| substring extraction | `s.substring(N/2, N)` | $c_9$ |
| string concatenation | `s + t` | $c_{10}\ N$ |

**Novice mistake.** Abusive string concatenation.

## Example: 1-sum

**Q.** How many instructions as a function of N?

```
int count = 0;
for (int i = 0; i < N; i++)
   if (a[i] == 0) count++;
```

| operation | frequency |
|---|---|
| variable declaration | 2 |
| assignment statement | 2 |
| less than compare | $N + 1$ |
| equal to compare | $N$ |
| array access | $N$ |
| increment | $\leq 2N$ |

between N (no zeros)
and 2N (all zeros)

## Example: 2-sum

Q. How many instructions as a function of N?

```
int count = 0;
for (int i = 0; i < N; i++)
   for (int j = i+1; j < N; j++)
      if (a[i] + a[j] == 0) count++;
```

| operation | frequency |
|---|---|
| variable declaration | $N + 2$ |
| assignment statement | $N + 2$ |
| less than compare | $1/2\ (N + 1)\ (N + 2)$ |
| equal to compare | $1/2\ N\ (N - 1)$ |
| array access | $N\ (N - 1)$ |
| increment | $\leq\ N^2$ |

$$0 + 1 + 2 + \ldots + (N - 1)\ =\ \frac{1}{2} N\ (N - 1)$$

$$=\ \binom{N}{2}$$

tedious to count exactly

## Tilde notation

- Estimate running time (or memory) as a function of input size $N$.
- Ignore lower order terms.
  - when $N$ is large, terms are negligible
  - when $N$ is small, we don't care

Ex 1.     $6\,N^3 + 20\,N + 16$         $\sim\ 6\,N^3$

Ex 2.     $6\,N^3 + 100\,N^{4/3} + 56$     $\sim\ 6\,N^3$

Ex 3.     $6\,N^3 + 17\,N^2\,\lg N + 7\,N$    $\sim\ 6\,N^3$

discard lower-order terms
(e.g., N = 1000: 6 billion vs. 169 million)

Technical definition.   $f(N) \sim g(N)$ means   $\displaystyle\lim_{N \to \infty} \frac{f(N)}{g(N)} = 1$

## Example: 2-sum

Q. How long will it take as a function of N?

```
int count = 0;
for (int i = 0; i < N; i++)
   for (int j = i+1; j < N; j++)
      if (a[i] + a[j] == 0) count++;
```

← "inner loop"

| operation | frequency | time per op | total time |
|---|---|---|---|
| variable declaration | $\sim N$ | $c_1$ | $\sim c_1 N$ |
| assignment statement | $\sim N$ | $c_2$ | $\sim c_2 N$ |
| less than comparison | $\sim 1/2 \, N^2$ | $c_3$ | $\sim c_3 N^2$ |
| equal to comparison | $\sim 1/2 \, N^2$ | | |
| array access | $\sim N^2$ | $c_4$ | $\sim c_4 N^2$ |
| increment | $\leq N^2$ | $c_5$ | $\leq c_5 N^2$ |
| total | | | $\sim c N^2$ |

depends on input data ↗

## Example: 3-sum

Q. How many instructions as a function of N?

```
int count = 0;                              ~ 1

for (int i = 0; i < N; i++)                 ~ N

    for (int j = i+1; j < N; j++)           ~ N² / 2

        for (int k = j+1; k < N; k++)

            if (a[i] + a[j] + a[k] == 0)    $\binom{N}{3} = \frac{N(N-1)(N-2)}{3!}$

                count++;                     $\sim \frac{1}{6}N^3$
```

"inner loop"

may be in inner loop, depends on input data

Remark. Focus on instructions in inner loop; ignore everything else!

## Bounding the sum by an integral trick

Q. How to estimate a discrete sum?

A1. Take COS 340.

A2. Replace the sum with an integral, and use calculus!

Ex 1. 1 + 2 + ... + N.
$$\sum_{i=1}^{N} i \ \sim \ \int_{x=1}^{N} x \, dx \ \sim \ \frac{1}{2} N^2$$

Ex 2. 1 + 1/2 + 1/3 + ... + 1/N.
$$\sum_{i=1}^{N} \frac{1}{i} \ \sim \ \int_{x=1}^{N} \frac{1}{x} dx \ = \ \ln N$$

Ex 3. 3-sum triple loop.
$$\sum_{i=1}^{N} \sum_{j=i}^{N} \sum_{k=j}^{N} 1 \ \sim \ \int_{x=1}^{N} \int_{y=x}^{N} \int_{z=y}^{N} dz \, dy \, dx \ \sim \ \frac{1}{6} N^3$$

## Mathematical models for running time

In principle, accurate mathematical models are available.

In practice,

- Formulas can be complicated.
- Advanced mathematics might be required.
- Exact models best left for experts.

costs (depend on machine, compiler)

$$T_N = c_1 A + c_2 B + c_3 C + c_4 D + c_5 E$$

A = variable declarations
B = assignment statements
C = compare
D = array access
E = increment

frequencies
(depend on algorithm, input)

Bottom line.  We use approximate models in this course:  $T_N \sim c\ N^3$.

# Common order-of-growth hypotheses

To determine order-of-growth:

- Assume a power law $T_N \sim a\,N^b$.
- Estimate exponent $b$ with doubling hypothesis.
- Validate with mathematical analysis.

Ex.  `ThreeSumDeluxe.java`

Food for precept.  How is it implemented?

| N | time (seconds) |
|---|---|
| 1,000 | 0.26 |
| 2,000 | 2.16 |
| 4,000 | 17.18 |
| 8,000 | 137.76 |

`ThreeSum.java`

| N | time (seconds) |
|---|---|
| 1,000 | 0.43 |
| 2,000 | 0.53 |
| 4,000 | 1.01 |
| 8,000 | 2.87 |
| 16,000 | 11.00 |
| 32,000 | 44.64 |
| 64,000 | 177.48 |

`ThreeSumDeluxe.java`

## Common order-of-growth hypotheses

**Good news.** the small set of functions

$$1, \ \log N, \ N, \ N \log N, \ N^2, \ N^3, \text{and } 2^N$$

suffices to describe order-of-growth of typical algorithms.



*Orders of growth (log-log plot)*

# Common order-of-growth hypotheses

| growth rate | name | typical code framework | description | example | T(2N) / T(N) |
|---|---|---|---|---|---|
| 1 | constant | `a = b + c;` | statement | add two numbers | 1 |
| log N | logarithmic | `while (N > 1)`<br>`{    N = N / 2;   ...    }` | divide in half | binary search | ~ 1 |
| N | linear | `for (int i = 0; i < N; i++)`<br>`{   ...          }` | loop | find the maximum | 2 |
| N log N | linearithmic | [see mergesort lecture] | divide and conquer | mergesort | ~ 2 |
| $N^2$ | quadratic | `for (int i = 0; i < N; i++)`<br>`  for (int j = 0; j < N; j++)`<br>`    {   ...          }` | double loop | check all pairs | 4 |
| $N^3$ | cubic | `for (int i = 0; i < N; i++)`<br>`  for (int j = 0; j < N; j++)`<br>`    for (int k = 0; k < N; k++)`<br>`      {   ...          }` | triple loop | check all triples | 8 |
| $2^N$ | exponential | [see combinatorial search lecture] | exhaustive search | check all possibilities | T(N) |

# Practical implications of order-of-growth

| growth rate | name | description | effect on a program that runs for a few seconds | |
|---|---|---|---|---|
| | | | time for 100x more data | size for 100x faster computer |
| 1 | constant | independent of input size | - | - |
| log N | logarithmic | nearly independent of input size | - | - |
| N | linear | optimal for N inputs | a few minutes | 100x |
| N log N | linearithmic | nearly optimal for N inputs | a few minutes | 100x |
| $N^2$ | quadratic | not practical for large problems | several hours | 10x |
| $N^3$ | cubic | not practical for medium problems | several weeks | 4-5x |
| $2^N$ | exponential | useful only for tiny problems | forever | 1x |

# Types of analyses

Best case.  Lower bound on cost.

- Determined by "easiest" input.
- Provides a goal for all inputs.

Worst case.  Upper bound on cost.

- Determined by "most difficult" input.
- Provides guarantee for all inputs.

Average case.  "Expected" cost.

- Need a model for "random" input.
- Provides a way to predict performance.

Ex 1. Array accesses for brute-force 3-sum.

- Best: ~ $\frac{1}{2}N^3$
- Average: ~ $\frac{1}{2}N^3$
- Worst:  ~ $\frac{1}{2}N^3$

Ex 2. Compares for insertion sort.

- Best (ascending order):  ~ N.
- Average (random order): ~ $\frac{1}{4} N^2$
- Worst (descending order):  ~ $\frac{1}{2}N^2$

(details in Lecture 4)

## Commonly-used notations

| notation | provides | example | shorthand for | used to |
|----------|----------|---------|---------------|---------|
| Tilde | leading term | $\sim 10\,N^2$ | $10\,N^2$ <br> $10\,N^2 + 22\,N \log N$ <br> $10\,N^2 + 2\,N + 37$ | provide approximate model |
| Big Theta | asymptotic growth rate | $\Theta(N^2)$ | $N^2$ <br> $9000\,N^2$ <br> $5\,N^2 + 22\,N \log N + 3N$ | classify algorithms |
| Big Oh | $\Theta(N^2)$ and smaller | $O(N^2)$ | $N^2$ <br> $100\,N$ <br> $22\,N \log N + 3\,N$ | develop upper bounds |
| Big Omega | $\Theta(N^2)$ and larger | $\Omega(N^2)$ | $9000\,N^2$ <br> $N^5$ <br> $N^3 + 22\,N \log N + 3\,N$ | develop lower bounds |

**Common mistake.** Interpreting big-Oh as an approximate model.

# Tilde notation vs. big-Oh notation

We use tilde notation whenever possible.

- Big-Oh notation suppresses leading constant.
- Big-Oh notation only provides upper bound (not lower bound).

## Typical memory requirements for primitive types in Java

Bit.  0 or 1.

Byte.  8 bits.

Megabyte (MB).  1 million bytes.

Gigabyte (GB).  1 billion bytes.

| type | bytes |
|---|---|
| boolean | 1 |
| byte | 1 |
| char | 2 |
| int | 4 |
| float | 4 |
| long | 8 |
| double | 8 |

# Typical memory requirements for arrays in Java

Array overhead.  16 bytes.

| type | bytes |
|---|---|
| `char[]` | $2N + 16$ |
| `int[]` | $4N + 16$ |
| `double[]` | $8N + 16$ |

one-dimensional arrays

| type | bytes |
|---|---|
| `char[][]` | $2N^2 + 20N + 16$ |
| `int[][]` | $4N^2 + 20N + 16$ |
| `double[][]` | $8N^2 + 20N + 16$ |

two-dimensional arrays

Ex.  An N-by-N array of doubles consumes $\sim 8N^2$ bytes of memory.

Typical memory requirements for objects in Java

Object overhead.  8 bytes.

Reference.  4 bytes.

Ex 1.  A `Complex` object consumes 24 bytes of memory.

```
public class Complex
{
    private double re;
    private double im;
    ...
}
```

8 bytes overhead for object

8 bytes

8 bytes

_____

24 bytes

*24 bytes*

| object overhead |
|---|
| re |
| im |

double *values*

# Typical memory requirements for objects in Java

Object overhead.  8 bytes.

Reference.  4 bytes.

Ex 2.  A virgin `string` of length N consumes ~ 2N bytes of memory.

```
public class String
{
    private int offset;
    private int count;
    private int hash;
    private char[] value;
    ...
}
```

8 bytes overhead for object

4 bytes

4 bytes

4 bytes

4 bytes for reference
(plus 2N + 16 bytes for array)

_____

2N + 40 bytes

| object overhead |
|:---:|
| value |
| offset |
| count |
| hash |

← reference

int
values

Example 1

Q. How much memory does `QuickUWPC` use as a function of $N$?

A.

```
public class QuickUWPC
{
    private int[] id;
    private int[] sz;

    public QuickUWPC(int N)
    {
        id = new int[N];
        sz = new int[N];
        for (int i = 0; i < N; i++) id[i] = i;
        for (int i = 0; i < N; i++) sz[i] = 1;
    }

    public boolean find(int p, int q)
    { ... }

    public void unite(int p, int q)
    { ... }
}
```

## Example 2

Q. How much memory does this code fragment use as a function of $N$?

A.

```
...
int N = Integer.parseInt(args[0]);
for (int i = 0; i < N; i++) {
    int[] a = new int[N];
    ...
}
```

Remark. Java automatically reclaims memory when it is no longer in use.

not always easy for Java to know

## Turning the crank:  summary

In principle, accurate mathematical models are available.

In practice,  approximate mathematical models are easily achieved.

Timing may be flawed?
- Limits on experiments insignificant compared to other sciences.

- Mathematics might be difficult?
- Only a few functions seem to turn up.
- Doubling hypothesis cancels complicated constants.

Actual data might not match input model?
- Need to understand input to effectively process it.
- Approach 1: design for the worst case.
- Approach 2: randomize, depend on probabilistic guarantee.

# 1.3  Stacks and Queues

‣ stacks
‣ dynamic resizing
‣ queues
‣ generics
‣ iterators
‣ applications

## Stacks and queues

**Fundamental data types.**

- Values: sets of objects
- Operations: insert, remove, test if empty.
- Intent is clear when we insert.
- Which item do we remove?

LIFO = "last in first out"

**Stack.** Remove the item most recently added.

**Analogy.** Cafeteria trays, Web surfing.

FIFO = "first in first out"

**Queue.** Remove the item least recently added.

**Analogy.** Registrar's line.

## Client, implementation, interface

**Separate interface and implementation.**

Ex:  stack, queue, priority queue, symbol table, union-find, ….

**Benefits.**

- Client can't know details of implementation  ⇒
  client has many implementation from which to choose.
- Implementation can't know details of client needs  ⇒
  many clients can re-use the same implementation.
- Design:  creates modular, reusable libraries.
- Performance:  use optimized implementation where it matters.

> Client:  program using operations defined in interface.
> Implementation:  actual code implementing operations.
> Interface:  description of data type, basic operations.

‣ **stacks**

# Stacks

## Stack operations.

- **push()**        Insert a new item onto stack.
- **pop()**         Remove and return the item most recently added.
- **isEmpty()**     Is the stack empty?

push

pop

```
public static void main(String[] args)
{
    StackOfStrings stack = new StackOfStrings();
    while (!StdIn.isEmpty())
    {
        String item = StdIn.readString();
        if (item.equals("-")) StdOut.print(stack.pop());
        else                  stack.push(item);
    }
}
```

```
% more tobe.txt
to be or not to - be - - that - - - is

% java StackOfStrings < tobe.txt
to be not that or be
```

# Stack pop: linked-list implementation



```
String item = first.item;
```

"of"

```
first = first.next;
```

```
return item;
```

"of"

## Stack push:  linked-list implementation

first

best ⟶ the ⟶ was ⟶ it

first   oldfirst

best ⟶ the ⟶ was ⟶ it

`Node oldfirst = first;`

first   oldfirst

best ⟶ the ⟶ was ⟶ it

`first = new Node();`

first   oldfirst

of ⟶ best ⟶ the ⟶ was ⟶ it

```
first.item = "of";
first.next = oldfirst;
```

## Stack:  linked-list implementation

```java
public class StackOfStrings
{
   private Node first = null;

   private class Node
   {
      String item;
      Node next;
   }

   public boolean isEmpty()
   {   return first == null;   }

   public void push(String item)
   {
      Node oldfirst = first;
      first = new Node();
      first.item = item;
      first.next = oldfirst;
   }

   public String pop()
   {
      if (isEmpty()) throw new RuntimeException();
      String item = first.item;
      first = first.next;
      return item;
   }
}
```

"inner class" ←

stack underflow ←

# Stack: linked-list trace

# Stack:  array implementation

Array implementation of a stack.

- Use array `s[]` to store `N` items on stack.
- `push()`: add new item at `s[N]`.
- `pop()`: remove item from `s[N-1]`.

| `s[]` | it | was | the | best | of | times | *null* | *null* | *null* | *null* |
|-------|----|-----|-----|------|-----|-------|--------|--------|--------|--------|
|       | 0  | 1   | 2   | 3    | 4   | 5     | 6      | 7      | 8      | 9      |

N                                         capacity = 10

# Stack:  array implementation

```
public class StackOfStrings
{
   private String[] s;                    a cheat
   private int N = 0;                   (stay tuned)


   public StackOfStrings(int capacity)
   {   s = new String[capacity];   }


   public boolean isEmpty()
   {   return N == 0;   }


   public void push(String item)
   {   s[N++] = item;   }


   public String pop()
   {   return s[--N];   }

}
```

decrement N;
then use to index into array

```
public String pop()
{
   String item = s[--N];
   s[N] = null;
   return item;

}
```

this version avoids "loitering"

garbage collector only reclaims memory
if no outstanding references

## Stack: dynamic array implementation

**Problem.** Requiring client to provide capacity does not implement API!

**Q.** How to grow and shrink array?

**First try.**
- `push()`: increase size of `s[]` by 1.
- `pop()`: decrease size of `s[]` by 1.

**Too expensive.**
- Need to copy all item to a new array.
- Inserting first N items takes time proportional to $1 + 2 + \ldots + N \sim N^2/2$.

infeasible for large N

**Goal.** Ensure that array resizing happens infrequently.

## Stack: dynamic array implementation

Q. How to grow array?

"repeated doubling"

A. If array is full, create a new array of twice the size, and copy items.

```
public StackOfStrings() {   s = new String[2];   }

public void push(String item)
{
    if (N == s.length) resize(2 * s.length);
    s[N++] = item;
}

private void resize(int capacity)
{
    String[] dup = new String[capacity];
    for (int i = 0; i < N; i++)
        dup[i] = s[i];
    s = dup;
}
```

$1 + 2 + 4 + \ldots + N/2 + N \sim 2N$

Consequence. Inserting first N items takes time proportional to N (not $N^2$).

## Stack:  dynamic array implementation

Q.  How to shrink array?

First try.

- `push()`:  double size of `s[]` when array is full.
- `pop()`:    halve size of `s[]` when array is half full.

Too expensive

- Consider push-pop-push-pop-… sequence when array is full.
- Takes time proportional to N per operation.

"thrashing"

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| N = 5 | it | was | the | best | of | *null* | *null* | *null* |
| N = 4 | it | was | the | best | | | |
| N = 5 | it | was | the | best | of | *null* | *null* | *null* |
| N = 4 | it | was | the | best | | | |

## Stack: dynamic array implementation

Q. How to shrink array?

Efficient solution.

- `push()`: double size of `s[]` when array is full.
- `pop()`: halve size of `s[]` when array is one-quarter full.

```java
public String pop()
{
    String item = s[--N];
    s[N] = null;
    if (N > 0 && N == s.length/4) resize(s.length / 2);
    return item;
}
```

Invariant. Array is always between 25% and 100% full.

# Stack:  dynamic array implementation trace

| StdIn | StdOut | N | a.length | a | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| | | 0 | 1 | *null* | | | | | | | |
| to | | 1 | 1 | to | | | | | | | |
| be | | 2 | 2 | to | be | | | | | | |
| or | | 3 | 4 | to | be | or | *null* | | | | |
| not | | 4 | 4 | to | be | or | not | | | | |
| to | | 5 | 8 | to | be | or | not | to | *null* | *null* | *null* |
| – | to | 4 | 8 | to | be | or | not | *null* | *null* | *null* | *null* |
| be | | 5 | 8 | to | be | or | not | be | *null* | *null* | *null* |
| – | be | 4 | 8 | to | be | or | not | *null* | *null* | *null* | *null* |
| – | not | 3 | 8 | to | be | or | *null* | *null* | *null* | *null* | *null* |
| that | | 4 | 8 | to | be | or | that | *null* | *null* | *null* | *null* |
| – | that | 3 | 8 | to | be | or | *null* | *null* | *null* | *null* | *null* |
| – | or | 2 | 4 | to | be | *null* | *null* | | | | |
| – | be | 1 | 2 | to | *null* | | | | | | |
| is | | 2 | 2 | to | is | | | | | | |

# Amortized analysis

Amortized analysis.  Average running time per operation over
a worst-case sequence of operations.

Proposition.  Starting from empty data structure, any sequence of M push and
pop ops takes time proportional to M.

*running time for doubling stack with N items*

|  | worst | best | amortized |
|---|---|---|---|
| construct | 1 | 1 | 1 |
| push | N | 1 | 1 |
| pop | N | 1 | 1 |

doubling or shrinking

Remark.  Recall, WQUPC used amortized bound.

## Stack implementations: memory usage

Linked list implementation.  ~ 16N bytes.

```
private class Node
{
    String item;
    Node next;
}
```

8 bytes overhead for object

4 bytes

4 bytes
_____
16 bytes per item

Doubling array.  Between ~ 4N (100% full) and ~ 16N (25% full).

```
public class DoublingStackOfStrings
{
    private String[] s;
    private int N = 0;
    …
}
```

4 bytes × array size

4 bytes

Remark.  Our analysis doesn't include the memory for the items themselves.

## Stack implementations:  dynamic array vs. linked List

Tradeoffs.  Can implement with either array or linked list;
client can use interchangeably.  Which is better?

Linked list.

• Every operation takes constant time in worst-case.
• Uses extra time and space to deal with the links.

Array.

• Every operation takes constant amortized time.
• Less wasted space.

- ▶ stacks
- ▶ dynamic resizing
- ▶ **queues**
- ▶ generics
- ▶ iterators
- ▶ applications

## Queues

### Queue operations.

- **enqueue()**        Insert a new item onto queue.

- **dequeue()**        Delete and return the item least recently added.

- **isEmpty()**        Is the queue empty?

```
public static void main(String[] args)
{
    QueueOfStrings q = new QueueOfStrings();
    while (!StdIn.isEmpty())
    {
        String item = StdIn.readString();
        if (item.equals("-")) StdOut.print(q.dequeue());
        else                  q.enqueue(item);
    }
}
```



(CNN.COM GRAPHIC)

```
% more tobe.txt
to be or not to - be - - that - - - is

% java QueueOfStrings < tobe.txt
to be or not to be
```

# Queue dequeue: linked list implementation



first

last

| it | was | the | best | of |

`String item = first.item;`

`"it"`

first

last

| was | the | best | of |

`first = first.next;`

first

last

| was | the | best | of |

`return item;`

`"it"`

# Queue enqueue:  linked list implementation



```
Node oldlast = last;
```

```
last = new Node();
last.item = "of";
last.next = null;
```

```
oldlast.next = last;
```

## Queue: linked list implementation

```java
public class QueueOfStrings
{
   private Node first, last;

   private class Node
   {  /* same as in StackOfStrings */  }

   public boolean isEmpty()
   {  return first == null;  }

   public void enqueue(String item)
   {
      Node oldlast = last;
      last = new Node();
      last.item = item;
      last.next = null;
      if (isEmpty()) first = last;
      else           oldlast.next = last;
   }

   public String dequeue()
   {
      String item = first.item;
      first       = first.next;
      if (isEmpty()) last = null;
      return item;
   }
}
```

# Queue: dynamic array implementation

Array implementation of a queue.

- Use array `q[]` to store items in queue.
- `enqueue()`: add new item at `q[tail]`.
- `dequeue()`: remove item from `q[head]`.
- Update `head` and `tail` modulo the `capacity`.
- Add repeated doubling and shrinking.

| `q[]` | *null* | *null* | **the** | **best** | **of** | **times** | *null* | *null* | *null* | *null* |
|-------|--------|--------|---------|----------|--------|-----------|--------|--------|--------|--------|
|       | 0      | 1      | 2       | 3        | 4      | 5         | 6      | 7      | 8      | 9      |

    **head**          **tail**       **capacity = 10**

## Parameterized stack

We implemented: `StackOfStrings.`

We also want: `StackOfURLs, StackOfCustomers, StackOfInts,` etc?

Attempt 1. Implement a separate stack class for each type.
- Rewriting code is tedious and error-prone.
- Maintaining cut-and-pasted code is tedious and error-prone.

@#$*! most reasonable approach until Java 1.5.

[hence, used in Algorithms in Java, 3rd edition]

## Parameterized stack

We implemented: `StackOfStrings.`

We also want: `StackOfURLs, StackOfCustomers, StackOfInts,` etc?

Attempt 2. Implement a stack with items of type `Object.`

- Casting is required in client.
- Casting is error-prone: run-time error if types mismatch.

```
StackOfObjects s = new StackOfObjects();
Apple  a = new Apple();
Orange b = new Orange();
s.push(a);
s.push(b);
a = (Apple) (s.pop());                          ⟵———— run-time error
```

## Parameterized stack

We implemented: `StackOfStrings.`

We also want: `StackOfURLs, StackOfCustomers, StackOfInts,` etc?


Attempt 3. Java generics.

- Avoid casting in both client and implementation.

- Discover type mismatch errors at compile-time instead of run-time.

type parameter

```
Stack<Apple> s = new Stack<Apple>();
Apple  a = new Apple();
Orange b = new Orange();
s.push(a);
s.push(b);
a = s.pop();
```

compile-time error

Guiding principles. Welcome compile-time errors; avoid run-time errors.

# Generic stack: linked list implementation

```
public class LinkedStackOfStrings
{
   private Node first = null;

   private class Node
   {
      String item;
      Node next;
   }

   public boolean isEmpty()
   {   return first == null;   }

   public void push(String item)
   {
      Node oldfirst = first;
      first = new Node();
      first.item = item;
      first.next = oldfirst;
   }

   public String pop()
   {
      String item = first.item;
      first = first.next;
      return item;
   }
}
```

```
public class Stack<Item>
{
   private Node first = null;

   private class Node
   {
      Item item;
      Node next;
   }

   public boolean isEmpty()
   {   return first == null;   }

   public void push(Item item)
   {
      Node oldfirst = first;
      first = new Node();
      first.item = item;
      first.next = oldfirst;
   }

   public Item pop()
   {
      Item item = first.item;
      first = first.next;
      return item;
   }
}
```

generic type name

# Generic stack: array implementation

```
public class ArrayStackOfStrings
{
   private String[] s;
   private int N = 0;

   public StackOfStrings(int capacity)
   {   s = new String[capacity];   }

   public boolean isEmpty()
   {   return N == 0;   }

   public void push(String item)
   {   s[N++] = item;   }

   public String pop()
   {   return s[--N];   }
}
```

```
public class ArrayStack<Item>
{
   private Item[] s;
   private int N = 0;

   public Stack(int capacity)
   {   s = new Item[capacity];   }

   public boolean isEmpty()
   {   return N == 0;   }

   public void push(Item item)
   {   s[N++] = item;   }

   public Item pop()
   {   return s[--N];   }
}
```

the way it should be

@#$*! generic array creation not allowed in Java

# Generic stack: array implementation

```
public class ArrayStackOfStrings
{
    private String[] s;
    private int N = 0;

    public StackOfStrings(int capacity)
    {   s = new String[capacity];   }

    public boolean isEmpty()
    {   return N == 0;   }

    public void push(String item)
    {   s[N++] = item;   }

    public String pop()
    {   return s[--N];   }
}
```

```
public class ArrayStack<Item>
{
    private Item[] s;
    private int N = 0;

    public Stack(int capacity)
    {   s = (Item[]) new Object[capacity]; }

    public boolean isEmpty()
    { return N == 0; }

    public void push(Item item)
    {   s[N++] = item;   }

    public Item pop()
    {   return s[--N];   }
}
```

the way it is

the ugly cast

Generic data types: autoboxing

Q. What to do about primitive types?

Wrapper type.
- Each primitive type has a wrapper object type.
- Ex: `Integer` is wrapper type for `int`.

Autoboxing. Automatic cast between a primitive type and its wrapper.

Syntactic sugar. Behind-the-scenes casting.

```
Stack<Integer> s = new Stack<Integer>();
s.push(17);         // s.push(new Integer(17));
int a = s.pop();    // int a = s.pop().intValue();
```

Bottom line. Client code can use generic stack for any type of data.

Q. What does the following program print?

```java
public class Autoboxing {

    public static void cmp(Integer a, Integer b) {
        if        (a <  b) StdOut.printf("%d <  %d\n", a, b);
        else if (a == b) StdOut.printf("%d == %d\n", a, b);
        else             StdOut.printf("%d >  %d\n", a, b);
    }

    public static void main(String[] args) {
        cmp(new Integer(42), new Integer(42));
        cmp(43, 43);
        cmp(142, 142);
    }
}
```

```
% java Autoboxing
42 >  42
43 == 43
142 >  142
```

Best practice.  Avoid using wrapper types whenever possible.

## Generics

**Caveat.** Java generics can be mystifying at times.

```java
public class Collections
{
   ...
   public static<T> void copy(List<? super T> dest, List<? extends T> src)
   {
      for (int i = 0; i < src.size(); i++)
         dest.set(i, src.get(i));
   }
}
```

mixing generics with inheritance

*Speed Up The Java Development Process*

Java Generics
*and Collections*

O'REILLY®

*Maurice Naftalin & Philip Wadler*

**This course.** Restrict attention to "pure generics."

avoid mixing generics with inheritance

- stacks
- dynamic resizing
- queues
- generics
- **iterators**
- applications

# Iteration

**Design challenge.** Support iteration over stack items by client, without revealing the internal representation of the stack.

|  | i |  |  | N |  |  |  |
|---|---|---|---|---|---|---|---|

|  | it | was | the | best | of | times | *null* | *null* | *null* | *null* |
|---|---|---|---|---|---|---|---|---|---|---|
| s[] | | | | | | | | | | |
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

first → current

| of | → | best | → | the | → | was | → | it | → | *null* |
|---|---|---|---|---|---|---|---|---|---|---|

**Java solution.** Make stack implement the `Iterable` interface.

## Iterators

Q.  What is an `Iterable` ?

A.  Has a method that returns an `Iterator`.

```
public interface Iterable<Item>
{
    Iterator<Item> iterator();

}
```

Q.  What is an `Iterator` ?

A.  Has methods `hasNext()` and `next()`.

```
public interface Iterator<Item>
{
    boolean hasNext();
    Item next();
    void remove();    ←——  optional; use
                            at your own risk
}
```

Q.  Why make data structures `Iterable` ?

A.  Java supports elegant client code.

"foreach" statement

```
for (String s : stack)
    StdOut.println(s);
```

equivalent code

```
Iterator<String> i = stack.iterator();
while (i.hasNext())
{
    String s = i.next();
    StdOut.println(s);
}
```

# Stack iterator:  linked list implementation

```java
import java.util.Iterator;

public class Stack<Item> implements Iterable<Item>
{
    ...

    public Iterator<Item> iterator() { return new ListIterator();  }

    private class ListIterator implements Iterator<Item>
    {
        private Node current = first;

        public boolean hasNext() {   return current != null;  }
        public void remove()    {   /* not supported */      }
        public Item next()
        {
            Item item = current.item;
            current   = current.next;
            return item;
        }
    }
}
```

**first**

**current**

of → best → the → was → it → *null*

40

## Stack iterator:  array implementation

```java
import java.util.Iterator;

public class Stack<Item> implements Iterable<Item>
{
    …

    public Iterator<Item> iterator() { return new ArrayIterator(); }

    private class ArrayIterator implements Iterator<Item>
    {
        private int i = N;

        public boolean hasNext() {  return i > 0;        }
        public void remove()     {  /* not supported */  }
        public Item next()       {  return s[--i];       }
    }
}
```

|  |  |  |  |  |  |  |  |  |  |
|---|---|---|---|---|---|---|---|---|---|
|  |  |  |  |  | | **i** |  | **N** |  |
| **s[]** | it | was | the | best | of | times | *null* | *null* | *null* | *null* |
|  | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

▸ stacks

▸ dynamic resizing

▸ queues

▸ generics

▸ iterators

▸ **applications**

## Java collections library

`java.util.List` API.

- `boolean isEmpty()`        Is the list empty?
- `int size()`        Return number of items on the list.
- `void add(Item item)`        Insert a new item to end of list.
- `void add(int index, Item item)`   Insert item at specified index.
- `Item get(int index)`        Return item at given index.
- `Item remove(int index)`        Return and delete item at given index.
- `Item set(int index Item item)`    Replace element at given index.
- `boolean contains(Item item)`      Does the list contain the item?
- `Iterator<Item> iterator()`      Return iterator.
- …

## Implementations.

- `java.util.ArrayList` implements API using an array.
- `java.util.LinkedList` implements API using a (doubly) linked list.

## Java collections library

**java.util.Stack.**

- Supports `push()`, `pop()`, `size()`, `isEmpty()`, and iteration.
- Also implements `java.util.List` interface from previous slide,
  e.g., `set()`, `get()`, and `contains()`.
- Bloated and poorly-designed API ⟹ don't use.

**java.util.Queue.**

- An interface, not an implementation of a queue.

Best practices. Use our implementations of `Stack` and `Queue` if you need a stack or a queue.

War story (from COS 226)

Generate random open sites in an N-by-N percolation system.

- Jenny: pick (i, j) at random; if closed, repeat.

  Takes ~ $c_1 N^2$ seconds.

- Kenny: maintain a `java.util.ArrayList` of open sites.

  Pick an index at random and delete.

  Takes ~ $c_1 N^4$ seconds.

Q. Why is Kenny's code so slow?

Lesson. Don't use a library until you understand its API!

COS 226. Can't use a library until we've implemented it in class.

## Stack applications

**Real world applications.**

- Parsing in a compiler.
- Java virtual machine.
- Undo in a word processor.
- Back button in a Web browser.
- PostScript language for printers.
- Implementing function calls in a compiler.

## Function calls

How a compiler implements a function.

- Function call: push local environment and return address.
- Return: pop return address and local environment.

Recursive function. Function that calls itself.

Note. Can always use an explicit stack to remove recursion.

p = 216, q = 192

gcd (216, 192)

```
static int gcd(int p, int q) {
    if (q == 0) return p;
    else
}
```

p = 192, q = 24

gcd (192, 24)

```
static int gcd(int p, int q) {
    if (q == 0) return p;
    else
}
```

p = 24, q = 0

gcd (24, 0)

```
static int gcd(int p, int q) {
    if (q == 0) return p;
    else return gcd(q, p % q);
}
```

# Arithmetic expression evaluation

**Goal.** Evaluate infix expressions.

```
( 1 + ( ( 2 + 3 ) * ( 4 * 5 ) ) )
```

operand      operator

**Two-stack algorithm.** [E. W. Dijkstra]

- Value: push onto the value stack.
- Operator: push onto the operator stack.
- Left parens: ignore.
- Right parens: pop operator and two values; push the result of applying that operator to those values onto the operand stack.

**Context.** An interpreter!

value stack
operator stack

| | |
|---|---|
| | ( 1 + ( ( 2 + 3 ) * ( 4 * 5 ) ) ) |
| 1 | + ( ( 2 + 3 ) * ( 4 * 5 ) ) ) |
| 1 / + | ( ( 2 + 3 ) * ( 4 * 5 ) ) ) |
| 1 2 / + | + 3 ) * ( 4 * 5 ) ) ) |
| 1 2 / + + | 3 ) * ( 4 * 5 ) ) ) |
| 1 2 3 / + + | ) * ( 4 * 5 ) ) ) |
| 1 5 / + | * ( 4 * 5 ) ) ) |
| 1 5 / + * | ( 4 * 5 ) ) ) |
| 1 5 4 / + * | * 5 ) ) ) |
| 1 5 4 / + * * | 5 ) ) ) |
| 1 5 4 5 / + * * | ) ) ) |
| 1 5 20 / + * | ) ) |
| 1 100 / + | ) |
| 101 | |

48

## Arithmetic expression evaluation

```java
public class Evaluate
{
   public static void main(String[] args)
   {
      Stack<String> ops  = new Stack<String>();
      Stack<Double> vals = new Stack<Double>();
      while (!StdIn.isEmpty()) {
         String s = StdIn.readString();
         if      (s.equals("("))                  ;
         else if (s.equals("+"))    ops.push(s);
         else if (s.equals("*"))    ops.push(s);
         else if (s.equals(")"))
         {
            String op = ops.pop();
            if      (op.equals("+")) vals.push(vals.pop() + vals.pop());
            else if (op.equals("*")) vals.push(vals.pop() * vals.pop());
         }
         else vals.push(Double.parseDouble(s));
      }
      StdOut.println(vals.pop());
   }
}
```

```
%  java Evaluate
( 1 + ( ( 2 + 3 ) * ( 4 * 5 ) ) )
101.0
```

## Correctness

Q. Why correct?

A. When algorithm encounters an operator surrounded by two values within parentheses, it leaves the result on the value stack.

```
( 1 + ( ( 2 + 3 ) * ( 4 * 5 ) ) )
```

as if the original input were:

```
( 1 + ( 5 * ( 4 * 5 ) ) )
```

Repeating the argument:

```
( 1 + ( 5 * 20 ) )
( 1 + 100 )
101
```

Extensions. More ops, precedence order, associativity.

## Stack-based programming languages

Observation 1. The 2-stack algorithm computes the same value if the operator occurs after the two values.

```
( 1 ( ( 2 3 + ) ( 4 5 * ) * ) + )
```

Observation 2. All of the parentheses are redundant!

```
1 2 3 + 4 5 * * +
```

Jan Lukasiewicz

Bottom line. Postfix or "reverse Polish" notation.

Applications. Postscript, Forth, calculators, Java virtual machine, …

## PostScript

Page description language.

- Explicit stack.
- Full computational model
- Graphics engine.

Basics.

- `%!`: "I am a PostScript program."
- Literal: "push me on the stack."
- Function calls take arguments from stack.
- Turtle graphics built in.

```
%!
72 72 moveto
0 72 rlineto
72 0 rlineto
0 -72 rlineto
-72 0 rlineto
2 setlinewidth
stroke
```

## PostScript

### Data types.

- Basic:  integer, floating point, boolean, ...
- Graphics:  font, path, curve, ....
- Full set of built-in operators.

### Text and strings.

- Full font support.
- `show` (display a string, using current font).

  `System.out.print()`

- `cvs` (convert anything to a string).

  `toString()`

```
%!
/Helvetica-Bold findfont 16 scalefont setfont
72 168 moveto
(Square root of 2:) show
72 144 moveto
2 sqrt 10 string cvs show
```

**Square root of 2:**
**1.41421**

53

# PostScript

## Variables (and functions).

- Identifiers start with /.
- **def** operator associates id with value.
- Braces.
- args on stack.

function definition →

```
%!
/box
{
  /sz exch def
  0 sz rlineto
  sz 0 rlineto
  0 sz neg rlineto
  sz neg 0 rlineto
} def

72 144 moveto
72 box
288 288 moveto
144 box
2 setlinewidth
stroke
```

function calls

54

## PostScript

### For loop.

- "from, increment, to" on stack.
- Loop body in braces.
- `for` operator.

### If-else conditional.

- Boolean on stack.
- Alternatives in braces.
- `if` operator.

... (hundreds of operators)

```
%!
\box
{
   ...
}

1 1 20
{ 19 mul dup 2 add moveto 72 box }
for
stroke
```

## PostScript

Application 1. All figures in Algorithms in Java, 3rd edition: figures created directly in PostScript.

```
%!
72 72 translate

/kochR
  {
    2 copy ge { dup 0 rlineto }
      {
        3 div
        2 copy kochR 60 rotate
        2 copy kochR -120 rotate
        2 copy kochR 60 rotate
        2 copy kochR
      } ifelse
    pop pop
  } def

0   0 moveto   81 243 kochR
0  81 moveto   27 243 kochR
0 162 moveto    9 243 kochR
0 243 moveto    1 243 kochR
stroke
```

See page 218

Application 2. All figures in Algorithms, 4th edition: enhanced version of `StdDraw` saves to PostScript for vector graphics.

## Queue applications

### Familiar applications.

- iTunes playlist.
- Data buffers (iPod, TiVo).
- Asynchronous data transfer (file IO, pipes, sockets).
- Dispensing requests on a shared resource (printer, processor).

### Simulations of the real world.

- Traffic analysis.
- Waiting times of customers at call center.
- Determining number of cashiers to have at a supermarket.

# M/M/1 queuing model

## M/M/1 queue.

- Customers arrive according to Poisson process at rate of $\lambda$ per minute.
- Customers are serviced with rate of $\mu$ per minute.

interarrival time has exponential distribution $\Pr[X \leq x] = 1 - e^{-\lambda x}$
service time has exponential distribution $\Pr[X \leq x] = 1 - e^{-\mu x}$

Arrival rate $\lambda$ → | | | | | | | → ⬤ → Departure rate $\mu$

Infinite queue       Server

Q. What is average wait time W of a customer in system?

Q. What is average number of customers L in system?

# M/M/1 queuing model:  example simulation



*time (seconds)*

|   | *arrival* | *departure* | *wait* |
|---|---|---|---|
| 0 | 0 | 5 | 5 |
| 1 | 2 | 10 | 8 |
| 2 | 7 | 15 | 8 |
| 3 | 17 | 23 | 6 |
| 4 | 19 | 28 | 9 |
| 5 | 21 | 30 | 9 |

## M/M/1 queuing model:  event-based simulation

```java
public class MM1Queue
{
    public static void main(String[] args) {
        double lambda = Double.parseDouble(args[0]);   // arrival rate
        double mu     = Double.parseDouble(args[1]);   // service rate
        double nextArrival = StdRandom.exp(lambda);
        double nextService = nextArrival + StdRandom.exp(mu);

        Queue<Double> queue = new Queue<Double>();
        Histogram hist = new Histogram("M/M/1 Queue", 60);

        while (true)
        {
            while (nextArrival < nextService)
            {
                queue.enqueue(nextArrival);
                nextArrival += StdRandom.exp(lambda);
            }

            double arrival = queue.dequeue();
            double wait = nextService - arrival;
            hist.addDataPoint(Math.min(60,  (int) (Math.round(wait))));
            if (queue.isEmpty()) nextService = nextArrival + StdRandom.exp(mu);
            else                 nextService = nextService + StdRandom.exp(mu);
        }
    }
}
```

next event is an arrival

next event is a service completion

# M/M/1 queuing model: experiments

Observation. If service rate $\mu$ is much larger than arrival rate $\lambda$, customers gets good service.



```
% java MM1Queue .2 .333
```

# M/M/1 queuing model: experiments

**Observation.** As service rate μ approaches arrival rate λ, services goes to h***.

```
% java MM1Queue .2 .25
```

# M/M/1 queuing model: experiments

Observation. As service rate μ approaches arrival rate λ, services goes to h***.

```
% java MM1Queue .2 .21
```

# M/M/1 queuing model:  analysis

M/M/1 queue.  Exact formulas known.

wait time W and queue length L approach infinity
as service rate approaches arrival rate

Little's Law

$$W \; = \; \frac{1}{\mu - \lambda} \; , \quad L \; = \; \lambda \, W$$



More complicated queueing models.  Event-based simulation essential!

Queueing theory.  See ORF 309.

# 2.1 Elementary Sorts

algorithms Sorting Write number array sort Time Insertion values order Selection Running use performance method code data Shellsort arrays one two exchanges see Example exercise input

- ▸ rules of the game
- ▸ selection sort
- ▸ insertion sort
- ▸ sorting challenges
- ▸ shellsort

## Sorting problem

Ex. Student record in a University.

| | | | | |
|---|---|---|---|---|
| Fox | 1 | A | 243-456-9091 | 101 Brown |
| Quilici | 1 | C | 343-987-5642 | 32 McCosh |
| Chen | 2 | A | 884-232-5341 | 11 Dickinson |
| Furia | 3 | A | 766-093-9873 | 22 Brown |
| Kanaga | 3 | B | 898-122-9643 | 343 Forbes |
| Andrews | 3 | A | 874-088-1212 | 121 Whitman |
| Rohde | 3 | A | 232-343-5555 | 115 Holder |
| Battle | 4 | C | 991-878-4944 | 308 Blair |
| Aaron | 4 | A | 664-480-0023 | 097 Little |
| Gazsi | 4 | B | 665-303-0266 | 113 Walker |

file →   record →   key →

Sort. Rearrange array of N objects into ascending order.

| | | | | |
|---|---|---|---|---|
| Aaron | 4 | A | 664-480-0023 | 097 Little |
| Andrews | 3 | A | 874-088-1212 | 121 Whitman |
| Battle | 4 | C | 991-878-4944 | 308 Blair |
| Chen | 2 | A | 884-232-5341 | 11 Dickinson |
| Fox | 1 | A | 243-456-9091 | 101 Brown |
| Furia | 3 | A | 766-093-9873 | 22 Brown |
| Gazsi | 4 | B | 665-303-0266 | 113 Walker |
| Kanaga | 3 | B | 898-122-9643 | 343 Forbes |
| Rohde | 3 | A | 232-343-5555 | 115 Holder |
| Quilici | 1 | C | 343-987-5642 | 32 McCosh |

## Sample sort client

Goal.  Sort any type of data.

Ex 1.  Sort random numbers in ascending order.

```
public class Experiment
{
   public static void main(String[] args)
   {
      int N = Integer.parseInt(args[0]);
      Double[] a = new Double[N];
      for (int i = 0; i < N; i++)
         a[i] = StdRandom.uniform();
      Insertion.sort(a);
      for (int i = 0; i < N; i++)
         StdOut.println(a[i]);
   }

}
```

```
% java Experiment 10
0.08614716385210452
0.09054270895414829
0.10708746304898642
0.21166190071646818
0.363292849257276
0.460954145685913
0.5340026311350087
0.7216129793703496
0.9003500354411443
0.9293994908845686
```

## Sample sort client

Goal.  Sort any type of data.

Ex 2.  Sort strings from standard input in alphabetical order.

```
public class StringSorter
{
   public static void main(String[] args)
   {
      String[] a = StdIn.readAll().split("\\s+");
      Insertion.sort(a);
      for (int i = 0; i < a.length; i++)
         StdOut.println(a[i]);
   }
}
```

```
% more words3.txt
bed bug dad yet zoo ... all bad yes

% java StringSorter < words.txt
all bad bed bug dad ... yes yet zoo
```

## Sample sort client

Goal. Sort any type of data.

Ex 3. Sort the files in a given directory by filename.

```java
import java.io.File;
public class FileSorter
{
    public static void main(String[] args)
    {
        File directory = new File(args[0]);
        File[] files = directory.listFiles();
        Insertion.sort(files);
        for (int i = 0; i < files.length; i++)
            StdOut.println(files[i].getName());
    }
}
```

```
% java FileSorter .
Insertion.class
Insertion.java
InsertionX.class
InsertionX.java
Selection.class
Selection.java
Shell.class
Shell.java
ShellX.class
ShellX.java
```

## Callbacks

Goal.  Sort any type of data.

Q.  How can sort know to compare data of type `String`, `Double`, and `File` without any information about the type of an item?

Callbacks.
- Client passes array of objects to sorting routine.
- Sorting routine calls back object's compare function as needed.

Implementing callbacks.
- Java:  interfaces.
- C:  function pointers.
- C++:  class-type functors.
- ML:  first-class functions and functors.

# Callbacks: roadmap

client

```
import java.io.File;
public class FileSorter
{
    public static void main(String[] args)
    {
        File directory = new File(args[0]);
        File[] files = directory.listFiles();
        Insertion.sort(files);
        for (int i = 0; i < files.length; i++)
            StdOut.println(files[i].getName());
    }
}
```

object implementation

```
public class File
implements Comparable<File>
{
    ...
    public int compareTo(File b)
    {
        ...
        return -1;
        ...
        return +1;
        ...
        return 0;
    }
}
```

built in to Java

interface

```
public interface Comparable<Item>
{
    public int compareTo(Item that);
}
```

sort implementation

```
public static void sort(Comparable[] a)
{
    int N = a.length;
    for (int i = 0; i < N; i++)
        for (int j = i; j > 0; j--)
            if (a[j].compareTo(a[j-1]) < 0)
                exch(a, j, j-1);
            else break;
}
```

key point: no reference to `File`

## Comparable interface API

**Comparable interface.** Implement `compareTo()` so that `v.compareTo(w)`:

- Returns a negative integer if `v` is less than `w`.
- Returns a positive integer if `v` is greater than `w`.
- Returns zero if `v` is equal to `w`.
- Throw an exception if incompatible types or either is `null`.

```
public interface Comparable<Item>
{   public int compareTo(Item that);   }
```

**Required properties.** Must ensure a total order.

- Reflexive: $(v = v)$.
- Antisymmetric: if $(v < w)$ then $(w > v)$; if $(v = w)$ then $(w = v)$.
- Transitive: if $(v \leq w)$ and $(w \leq x)$ then $(v \leq x)$.

**Built-in comparable types.** `String`, `Double`, `Integer`, `Date`, `File`, ...

**User-defined comparable types.** Implement the `Comparable` interface.

# Implementing the Comparable interface: example 1

Date data type. Simplified version of `java.util.Date`.

```java
public class Date implements Comparable<Date>
{
   private final int month, day, year;

   public Date(int m, int d, int y)
   {
      month = m;
      day   = d;
      year  = y;
   }

   public int compareTo(Date that)
   {
      if (this.year  < that.year ) return -1;
      if (this.year  > that.year ) return +1;
      if (this.month < that.month) return -1;
      if (this.month > that.month) return +1;
      if (this.day   < that.day  ) return -1;
      if (this.day   > that.day  ) return +1;
      return 0;
   }
}
```

only compare dates to other dates

# Implementing the Comparable interface:  example 2

## Domain names.

- Subdomain: `bolle.cs.princeton.edu.`
- Reverse subdomain:  `edu.princeton.cs.bolle.`
- Sort by reverse subdomain to group by category.

```
public class Domain implements Comparable<Domain>
{
   private final String[] fields;
   private final int N;

   public Domain(String name)
   {
      fields = name.split("\\.");
      N = fields.length;
   }

   public int compareTo(Domain that)
   {
      for (int i = 0; i < Math.min(this.N, that.N); i++)
      {
         String s = fields[this.N - i - 1];
         String t = fields[that.N - i - 1];
         int cmp = s.compareTo(t);
         if      (cmp < 0) return -1;
         else if (cmp > 0) return +1;
      }
      return this.N - that.N;
   }
}
```

only use this trick when no danger of overflow

*subdomains*

```
ee.princeton.edu
cs.princeton.edu
princeton.edu
cnn.com
google.com
apple.com
www.cs.princeton.edu
bolle.cs.princeton.edu
```

*reverse-sorted subdomains*

```
com.apple
com.cnn
com.google
edu.princeton
edu.princeton.cs
edu.princeton.cs.bolle
edu.princeton.cs.www
edu.princeton.ee
```

# Two useful sorting abstractions

Helper functions.  Refer to data through compares and exchanges.

Less.  Is object `v` less than `w` ?

```
private static boolean less(Comparable v, Comparable w)
{   return v.compareTo(w) < 0;   }
```

Exchange.  Swap object in array `a[]` at index `i` with the one at index `j`.

```
private static void exch(Comparable[] a, int i, int j)
{
    Comparable t = a[i];
    a[i] = a[j];
    a[j] = t;
}
```

## Testing

Q. How to test if an array is sorted?

```
private static boolean isSorted(Comparable[] a)
{
    for (int i = 1; i < a.length; i++)
        if (less(a[i], a[i-1])) return false;
    return true;
}
```

Q. If the sorting algorithm passes the test, did it correctly sort its input?

A. Yes, if data accessed only through `exch()` and `less()`.

13

## Selection sort inner loop

To maintain algorithm invariants:

- Move the pointer to the right.

```
i++;
```



in final order

- Identify index of minimum item on right.

```
int min = i;
for (int j = i+1; j < N; j++)
    if (less(a[j], a[min]))
        min = j;
```



in final order

- Exchange into position.

```
exch(a, i, min);
```



in final order

## Selection sort:  Java implementation

```java
public class Selection {

   public static void sort(Comparable[] a)
   {
      int N = a.length;
      for (int i = 0; i < N; i++)
      {
         int min = i;
         for (int j = i+1; j < N; j++)
            if (less(a[j], a[min]))
               min = j;
         exch(a, i, min);
      }
   }

   private static boolean less(Comparable v, Comparable w)
   {  /* as before */  }

   private static void exch(Comparable[] a, int i, int j)
   {  /* as before */  }
}
```

## Selection sort: mathematical analysis

**Proposition A.** Selection sort uses $(N-1) + (N-2) + ... + 1 + 0 \sim N^2/2$ compares and N exchanges.



| i | min | a[] |   |   |   |   |   |   |   |   |   |   |
|----|-----|---|---|---|---|---|---|---|---|---|---|----|
|    |     | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|    |     | S | O | R | T | E | X | A | M | P | L | E |
| 0  | 6   | S | O | R | T | E | X | A | M | P | L | E |
| 1  | 4   | A | O | R | T | E | X | S | M | P | L | E |
| 2  | 10  | A | E | R | T | O | X | S | M | P | L | E |
| 3  | 9   | A | E | E | T | O | X | S | M | P | L | R |
| 4  | 7   | A | E | E | L | O | X | S | M | P | T | R |
| 5  | 7   | A | E | E | L | M | X | S | O | P | T | R |
| 6  | 8   | A | E | E | L | M | O | S | X | P | T | R |
| 7  | 10  | A | E | E | L | M | O | P | X | S | T | R |
| 8  | 8   | A | E | E | L | M | O | P | R | S | T | X |
| 9  | 9   | A | E | E | L | M | O | P | R | S | T | X |
| 10 | 10  | A | E | E | L | M | O | P | R | S | T | X |
|    |     | A | E | E | L | M | O | P | R | S | T | X |

*entries in black are examined to find the minimum*

*entries in red are a[min]*

*entries in gray are in final position*

**Trace of selection sort (array contents just after each exchange)**

Running time insensitive to input. Quadratic time, even if array is presorted.

Data movement is minimal. Linear number of exchanges.

# Selection sort animations

20 random elements



▲ algorithm position

in final order

not in final order

http://www.sorting-algorithms.com/selection-sort

# Selection sort animations

20 partially-sorted elements



▲ algorithm position

━━━ in final order

━━━ not in final order

http://www.sorting-algorithms.com/selection-sort

## Insertion sort

**Algorithm.** ↑ scans from left to right.

**Invariants.**

- Elements to the left of ↑ (including ↑) are in ascending order.
- Elements to the right of ↑ have not yet been seen.



in order          ↑                    not yet seen

# Insertion sort inner loop

To maintain algorithm invariants:

- Move the pointer to the right.

  ```
  i++;
  ```



in order        not yet seen

- Moving from right to left, exchange
  `a[i]` with each larger element to its left.

  ```
  for (int j = i; j > 0; j--)
     if (less(a[j], a[j-1]))
        exch(a, j, j-1);
     else break;
  ```



in order        not yet seen

## Insertion sort: Java implementation

```java
public class Insertion {

    public static void sort(Comparable[] a)
    {
        int N = a.length;
        for (int i = 0; i < N; i++)
            for (int j = i; j > 0; j--)
                if (less(a[j], a[j-1]))
                    exch(a, j, j-1);
                else break;
    }

    private static boolean less(Comparable v, Comparable w)
    {  /* as before */  }

    private static void exch(Comparable[] a, int i, int j)
    {  /* as before */  }
}
```

# Insertion sort:  mathematical analysis

**Proposition B.**  To sort a randomly-ordered array with distinct keys, insertion sort uses ~ $N^2/4$ compares and $N^2/4$ exchanges on average.

**Pf.**  For randomly-ordered data, we expect each element to move halfway back.

|  i |  j | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|----|----|---|---|---|---|---|---|---|---|---|---|----|
|    |    | S | O | R | T | E | X | A | M | P | L | E |
|  1 |  0 | O | S | R | T | E | X | A | M | P | L | E |
|  2 |  1 | O | R | S | T | E | X | A | M | P | L | E |
|  3 |  3 | O | R | S | T | E | X | A | M | P | L | E |
|  4 |  0 | E | O | R | S | T | X | A | M | P | L | E |
|  5 |  5 | E | O | R | S | T | X | A | M | P | L | E |
|  6 |  0 | A | E | O | R | S | T | X | M | P | L | E |
|  7 |  2 | A | E | M | O | R | S | T | X | P | L | E |
|  8 |  4 | A | E | M | O | P | R | S | T | X | L | E |
|  9 |  2 | A | E | L | M | O | P | R | S | T | X | E |
| 10 |  2 | A | E | E | L | M | O | P | R | S | T | X |
|    |    | A | E | E | L | M | O | P | R | S | T | X |

a[]

*entries in gray do not move*

*entry in red is a[j]*

*entries in black moved one position right for insertion*

**Trace of insertion sort (array contents just after each insertion)**

# Insertion sort:  trace

a[ ]

| i | j | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  |  | A | S | O | M | E | W | H | A | T | L | O | N | G | E | R | I | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 0 | 0 | A | S | O | M | E | W | H | A | T | L | O | N | G | E | R | I | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 1 | 1 | A | S | O | M | E | W | H | A | T | L | O | N | G | E | R | I | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 2 | 1 | A | O | S | M | E | W | H | A | T | L | O | N | G | E | R | I | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 3 | 1 | A | M | O | S | E | W | H | A | T | L | O | N | G | E | R | I | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 4 | 1 | A | E | M | O | S | W | H | A | T | L | O | N | G | E | R | I | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 5 | 5 | A | E | M | O | S | W | H | A | T | L | O | N | G | E | R | I | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 6 | 2 | A | E | H | M | O | S | W | A | T | L | O | N | G | E | R | I | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 7 | 1 | A | A | E | H | M | O | S | W | T | L | O | N | G | E | R | I | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 8 | 7 | A | A | E | H | M | O | S | T | W | L | O | N | G | E | R | I | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 9 | 4 | A | A | E | H | L | M | O | S | T | W | O | N | G | E | R | I | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 10 | 7 | A | A | E | H | L | M | O | O | S | T | W | N | G | E | R | I | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 11 | 6 | A | A | E | H | L | M | N | O | O | S | T | W | G | E | R | I | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 12 | 3 | A | A | E | G | H | L | M | N | O | O | S | T | W | E | R | I | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 13 | 3 | A | A | E | E | G | H | L | M | N | O | O | S | T | W | R | I | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 14 | 11 | A | A | E | E | G | H | L | M | N | O | O | R | S | T | W | I | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 15 | 6 | A | A | E | E | G | H | I | L | M | N | O | O | R | S | T | W | N | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 16 | 10 | A | A | E | E | G | H | I | L | M | N | N | O | O | R | S | T | W | S | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 17 | 15 | A | A | E | E | G | H | I | L | M | N | N | O | O | R | S | S | T | W | E | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 18 | 4 | A | A | E | E | E | G | H | I | L | M | N | N | O | O | R | S | S | T | W | R | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 19 | 15 | A | A | E | E | E | G | H | I | L | M | N | N | O | O | R | R | S | S | T | W | T | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 20 | 19 | A | A | E | E | E | G | H | I | L | M | N | N | O | O | R | R | S | S | T | T | W | I | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 21 | 8 | A | A | E | E | E | G | H | I | I | L | M | N | N | O | O | R | R | S | S | T | T | W | O | N | S | O | R | T | E | X | A | M | P | L | E |
| 22 | 15 | A | A | E | E | E | G | H | I | I | L | M | N | N | O | O | O | R | R | S | S | T | T | W | N | S | O | R | T | E | X | A | M | P | L | E |
| 23 | 13 | A | A | E | E | E | G | H | I | I | L | M | N | N | N | O | O | O | R | R | S | S | T | T | W | S | O | R | T | E | X | A | M | P | L | E |
| 24 | 21 | A | A | E | E | E | G | H | I | I | L | M | N | N | N | O | O | O | R | R | S | S | S | T | T | W | O | R | T | E | X | A | M | P | L | E |
| 25 | 17 | A | A | E | E | E | G | H | I | I | L | M | N | N | N | O | O | O | O | R | R | S | S | S | T | T | W | R | T | E | X | A | M | P | L | E |
| 26 | 20 | A | A | E | E | E | G | H | I | I | L | M | N | N | N | O | O | O | O | R | R | R | S | S | S | T | T | W | T | E | X | A | M | P | L | E |
| 27 | 26 | A | A | E | E | E | G | H | I | I | L | M | N | N | N | O | O | O | O | R | R | R | S | S | S | T | T | T | W | E | X | A | M | P | L | E |
| 28 | 5 | A | A | E | E | E | E | G | H | I | I | L | M | N | N | N | O | O | O | O | R | R | R | S | S | S | T | T | T | W | X | A | M | P | L | E |
| 29 | 29 | A | A | E | E | E | E | G | H | I | I | L | M | N | N | N | O | O | O | O | R | R | R | S | S | S | T | T | T | W | X | A | M | P | L | E |
| 30 | 2 | A | A | A | E | E | E | E | G | H | I | I | L | M | N | N | N | O | O | O | O | R | R | R | S | S | S | T | T | T | W | X | M | P | L | E |
| 31 | 13 | A | A | A | E | E | E | E | G | H | I | I | L | M | M | N | N | N | O | O | O | O | R | R | R | S | S | S | T | T | T | W | X | P | L | E |
| 32 | 21 | A | A | A | E | E | E | E | G | H | I | I | L | M | M | N | N | N | O | O | O | O | P | R | R | R | S | S | S | T | T | T | W | X | L | E |
| 33 | 12 | A | A | A | E | E | E | E | G | H | I | I | L | L | M | M | N | N | N | O | O | O | O | P | R | R | R | S | S | S | T | T | T | W | X | E |
| 34 | 7 | A | A | A | E | E | E | E | E | G | H | I | I | L | L | M | M | N | N | N | O | O | O | O | P | R | R | R | S | S | S | T | T | T | W | X |
|  |  | A | A | A | E | E | E | E | E | G | H | I | I | L | L | M | M | N | N | N | O | O | O | O | P | R | R | R | S | S | S | T | T | T | W | X |

25

# Insertion sort animation

40 random elements

▲ algorithm position

in order

not yet seen

Insertion sort: best and worst case

Best case.  If the input is in ascending order, insertion sort makes
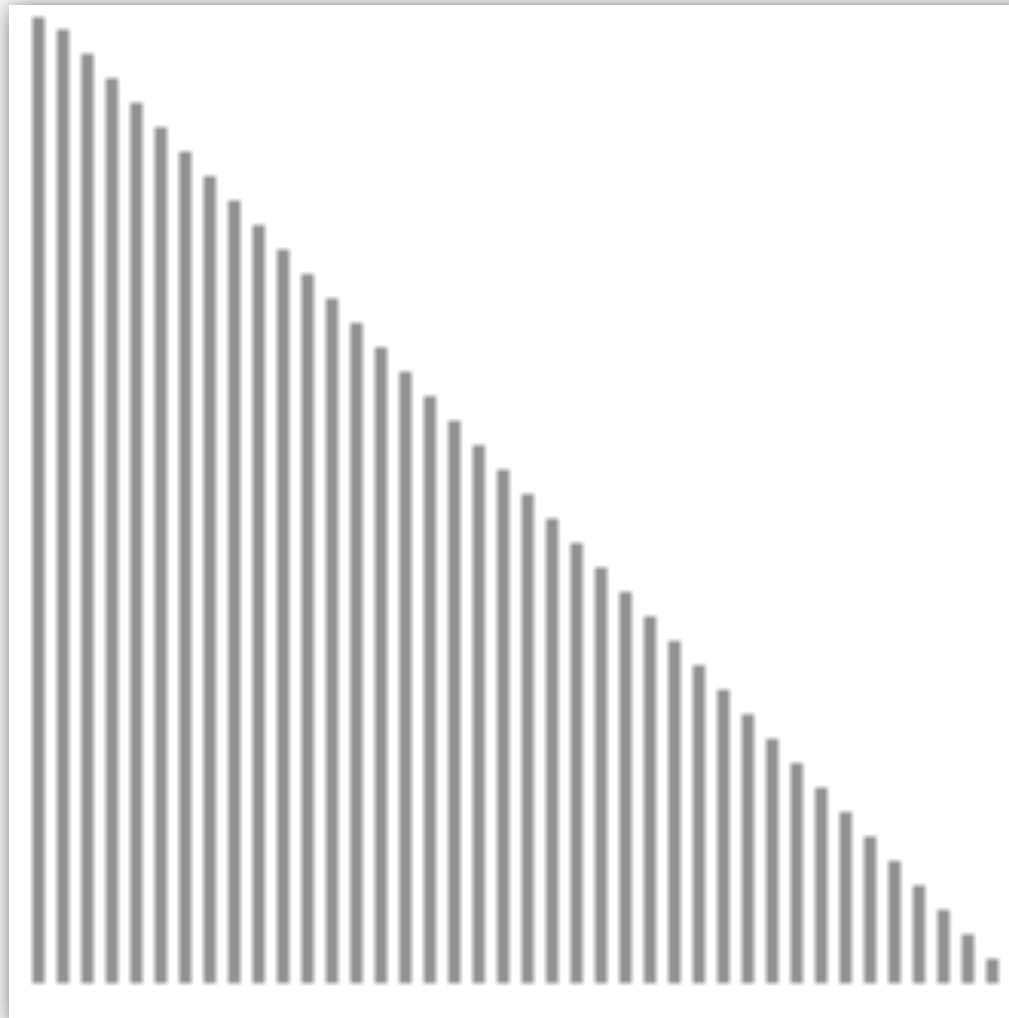N-1 compares and 0 exchanges.

```
A E E L M O P R S T X
```

Worst case.  If the input is in descending order (and no duplicates),
insertion sort makes ~ $N^2/2$ compares and ~ $N^2/2$ exchanges.

```
X T S R P O M L E E A
```

# Insertion sort animation

40 reverse-sorted elements



http://www.sorting-algorithms.com/insertion-sort

▲ algorithm position

in order

not yet seen

## Insertion sort: partially sorted inputs

Def.  An inversion is a pair of keys that are out of order.

A E E L M O T R X P S

T–R  T–P  T–S  R–P  X–P  X–S

(6 inversions)

Def. An array is partially sorted if the number of inversions is O(N).
- Ex 1.  A small array appended to a large sorted array.
- Ex 2. An array with only a few elements out of place.

Proposition C.  For partially-sorted arrays, insertion sort runs in linear time.

Pf.  Number of exchanges equals the number of inversions.

↑

number of compares = exchanges + (N-1)

# Insertion sort animation

40 partially-sorted elements



algorithm position
in order
not yet seen

http://www.sorting-algorithms.com/insertion-sort

# Sorting challenge 0

Input. Array of doubles.

Plot. Data proportional to length.

Name the sorting method.

- Insertion sort.
- Selection sort.



*gray entries are untouched*

*black entries are involved in compares*

# Sorting challenge 1

**Problem.** Sort a file of huge records with tiny keys.

**Ex.** Reorganize your MP3 files.

**Which sorting method to use?**

- System sort.
- Insertion sort.
- Selection sort.



| | | | | |
|---|---|---|---|---|
| Fox | 1 | A | 243-456-9091 | 101 Brown |
| Quilici | 1 | C | 343-987-5642 | 32 McCosh |
| Chen | 2 | A | 884-232-5341 | 11 Dickinson |
| Furia | 3 | A | 766-093-9873 | 22 Brown |
| Kanaga | 3 | B | 898-122-9643 | 343 Forbes |
| Andrews | 3 | A | 874-088-1212 | 121 Whitman |
| Rohde | 3 | A | 232-343-5555 | 115 Holder |
| Battle | 4 | C | 991-878-4944 | 308 Blair |
| Aaron | 4 | A | 664-480-0023 | 097 Little |
| Gazsi | 4 | B | 665-303-0266 | 113 Walker |

file → record → key →

## Sorting challenge 2

Problem.  Sort a huge randomly-ordered file of small records.

Ex.  Process transaction records for a phone company.

Which sorting method to use?

- System sort.

- Insertion sort.

- Selection sort.

## Sorting challenge 3

Problem.  Sort a huge number of tiny files (each file is independent).

Ex.  Daily customer transaction records.

Which sorting method to use?

- System sort.
- Insertion sort.
- Selection sort.

| | | | | |
|---|---|---|---|---|
| Fox | 1 | A | 243-456-9091 | 101 Brown |
| Quilici | 1 | C | 343-987-5642 | 32 McCosh |
| Chen | 2 | A | 884-232-5341 | 11 Dickinson |
| Furia | 3 | A | 766-093-9873 | 22 Brown |
| Kanaga | 3 | B | 898-122-9643 | 343 Forbes |
| Andrews | 3 | A | 874-088-1212 | 121 Whitman |
| Rohde | 3 | A | 232-343-5555 | 115 Holder |
| Battle | 4 | C | 991-878-4944 | 308 Blair |
| Aaron | 4 | A | 664-480-0023 | 097 Little |
| Gazsi | 4 | B | 665-303-0266 | 113 Walker |

file ➡

record ➡

key ➡

## Sorting challenge 4

Problem.  Sort a huge file that is already almost in order.

Ex.  Resort a huge database after a few changes.

Which sorting method to use?

- System sort.
- Insertion sort.
- Selection sort.



| file → | Fox | 1 | A | 243-456-9091 | 101 Brown |
|---|---|---|---|---|---|
| | Quilici | 1 | C | 343-987-5642 | 32 McCosh |
| | Chen | 2 | A | 884-232-5341 | 11 Dickinson |
| | Furia | 3 | A | 766-093-9873 | 22 Brown |
| | Kanaga | 3 | B | 898-122-9643 | 343 Forbes |
| record → | Andrews | 3 | A | 874-088-1212 | 121 Whitman |
| | Rohde | 3 | A | 232-343-5555 | 115 Holder |
| | Battle | 4 | C | 991-878-4944 | 308 Blair |
| key → | Aaron | 4 | A | 664-480-0023 | 097 Little |
| | Gazsi | 4 | B | 665-303-0266 | 113 Walker |

## Shellsort overview

**Idea.** Move elements more than one position at a time by h-sorting the array.

an h-sorted array is h interleaved sorted subsequences

h = 4

```
L   E   E   A   M   H   L   E   P   S   O   L   T   S   X   R
L───────────────M───────────────P───────────────T
    E───────────────H───────────────S───────────────S
        E───────────────L───────────────O───────────────X
            A───────────────E───────────────L───────────────R
```

**Shellsort.** h-sort the array for a decreasing sequence of values of h.

| | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| input | S | H | E | L | L | S | O | R | T | E | X | A | M | P | L | E |
| 13-sort | P | H | E | L | L | S | O | R | T | E | X | A | M | S | L | E |
| 4-sort | L | E | E | A | M | H | L | E | P | S | O | L | T | S | X | R |
| 1-sort | A | E | E | E | H | L | L | L | M | O | P | R | S | S | T | X |

38

# h-sorting

How to h-sort an array?  Insertion sort, with stride length h.

3-sorting an array

| M | O | L | E | E | X | A | S | P | R | T |
|---|---|---|---|---|---|---|---|---|---|---|
| E | O | L | M | E | X | A | S | P | R | T |
| E | E | L | M | O | X | A | S | P | R | T |
| E | E | L | M | O | X | A | S | P | R | T |
| A | E | L | E | O | X | M | S | P | R | T |
| A | E | L | E | O | X | M | S | P | R | T |
| A | E | L | E | O | P | M | S | X | R | T |
| A | E | L | E | O | P | M | S | X | R | T |
| A | E | L | E | O | P | M | S | X | R | T |
| A | E | L | E | O | P | M | S | X | R | T |

Why insertion sort?

• Big increments ⇒ small subarray.

• Small increments ⇒ nearly in order.  [stay tuned]

# Shellsort example: increments 7, 3, 1

**input**

| S | O | R | T | E | X | A | M | P | L | E |
|---|---|---|---|---|---|---|---|---|---|---|

**7-sort**

| S | O | R | T | E | X | A | M | P | L | E |
|---|---|---|---|---|---|---|---|---|---|---|
| M | O | R | T | E | X | A | S | P | L | E |
| M | O | R | T | E | X | A | S | P | L | E |
| M | O | L | T | E | X | A | S | P | R | E |
| M | O | L | E | E | X | A | S | P | R | T |

**3-sort**

| M | O | L | E | E | X | A | S | P | R | T |
|---|---|---|---|---|---|---|---|---|---|---|
| E | O | L | M | E | X | A | S | P | R | T |
| E | E | L | M | O | X | A | S | P | R | T |
| E | E | L | M | O | X | A | S | P | R | T |
| A | E | L | E | O | X | M | S | P | R | T |
| A | E | L | E | O | X | M | S | P | R | T |
| A | E | L | E | O | P | M | S | X | R | T |
| A | E | L | E | O | P | M | S | X | R | T |
| A | E | L | E | O | P | M | S | X | R | T |

**1-sort**

| A | E | L | E | O | P | M | S | X | R | T |
|---|---|---|---|---|---|---|---|---|---|---|
| A | E | L | E | O | P | M | S | X | R | T |
| A | E | L | E | O | P | M | S | X | R | T |
| A | E | E | L | O | P | M | S | X | R | T |
| A | E | E | L | O | P | M | S | X | R | T |
| A | E | E | L | O | P | M | S | X | R | T |
| A | E | E | L | M | O | P | S | X | R | T |
| A | E | E | L | M | O | P | S | X | R | T |
| A | E | E | L | M | O | P | S | X | R | T |
| A | E | E | L | M | O | P | R | S | X | T |
| A | E | E | L | M | O | P | R | S | T | X |

**result**

| A | E | E | L | M | O | P | R | S | T | X |
|---|---|---|---|---|---|---|---|---|---|---|

40

# Shellsort: intuition

**Proposition.** A g-sorted array remains g-sorted after h-sorting it.

**Pf.** Harder than you'd think!

7-sort

| M | O | R | T | E | X | A | S | P | L | E |
|---|---|---|---|---|---|---|---|---|---|---|
| M | O | R | T | E | X | A | S | P | L | E |
| M | O | L | T | E | X | A | S | P | R | E |
| M | O | L | E | E | X | A | S | P | R | T |
| M | O | L | E | E | X | A | S | P | R | T |

3-sort

| M | O | L | E | E | X | A | S | P | R | T |
|---|---|---|---|---|---|---|---|---|---|---|
| E | O | L | M | E | X | A | S | P | R | T |
| E | E | L | M | O | X | A | S | P | R | T |
| E | E | L | M | O | X | A | S | P | R | T |
| A | E | L | E | O | X | M | S | P | R | T |
| A | E | L | E | O | X | M | S | P | R | T |
| A | E | L | E | O | P | M | S | X | R | T |
| A | E | L | E | O | P | M | S | X | R | T |
| A | E | L | E | O | P | M | S | X | R | T |
| A | E | L | E | O | P | M | S | X | R | T |

still 7-sorted

What increments to use?

1, 2, 4, 8, 16, 32 . . .
No.

1, 3, 7, 15, 31, 63, . . .
Maybe.

→ 1, 4, 13, 40, 121, 364, . . .
OK, easy to compute 3x+1 sequence.

1, 5, 19, 41, 109, 209, 505, . . .
Tough to beat in empirical studies.

Interested in learning more?
- See Algs 3 section 6.8 or Knuth volume 3 for details.
- Consider doing a JP on the topic.

## Shellsort: Java implementation

```java
public class Shell
{
   public static void sort(Comparable[] a)
   {
      int N = a.length;

      int h = 1;
      while (h < N/3) h = 3*h + 1;  // 1, 4, 13, 40, 121, 364, 1093, ...

      while (h >= 1)
      {  // h-sort the array.
         for (int i = h; i < N; i++)
         {
            for (int j = i; j >= h && less(a[j], a[j-h]); j -= h)
               exch(a, j, j-h);
         }

         h = h/3;
      }
   }

   private static boolean less(Comparable v, Comparable w)
   {  /* as before */  }
   private static boolean void(Comparable[] a, int i, int j)
   {  /* as before */  }
}
```
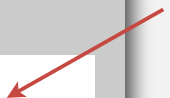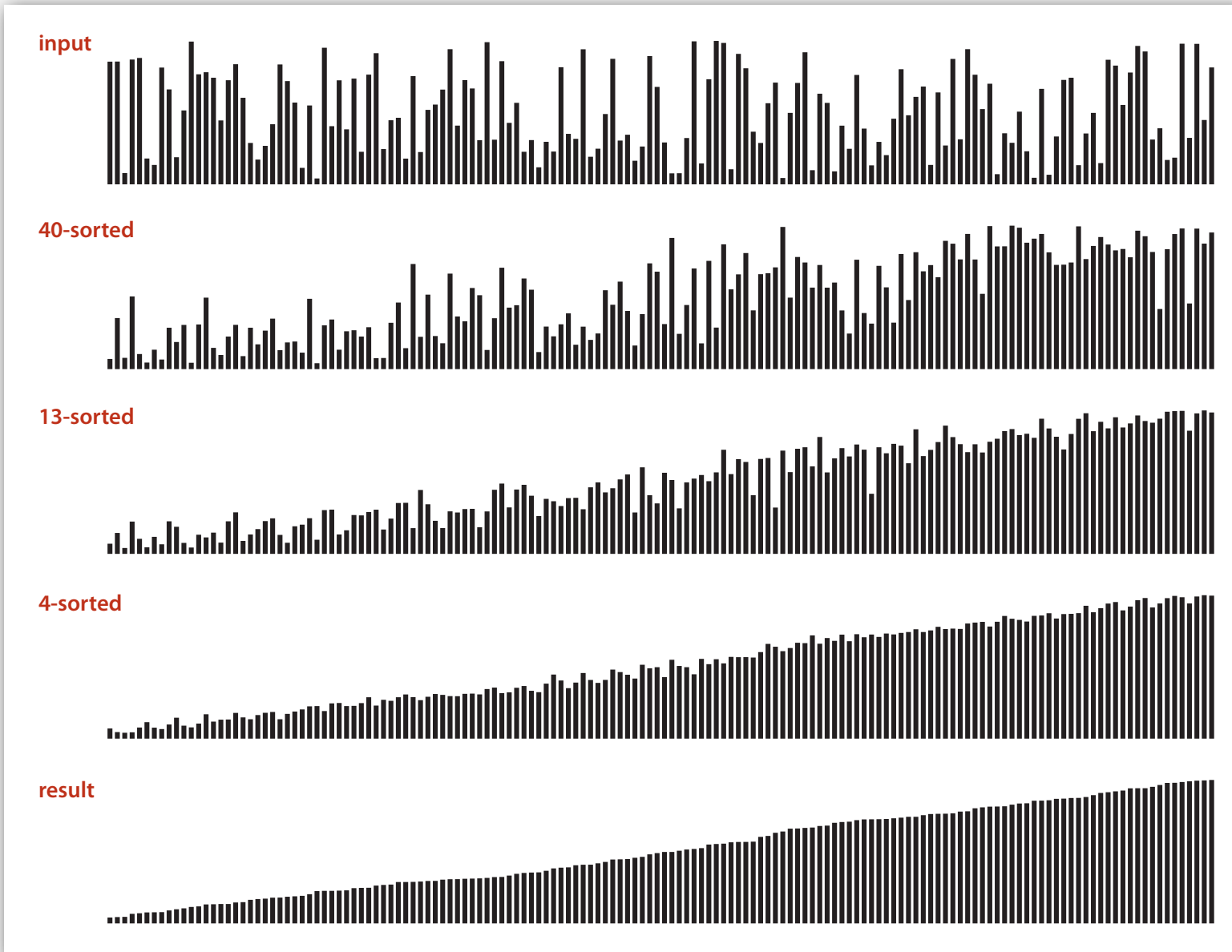
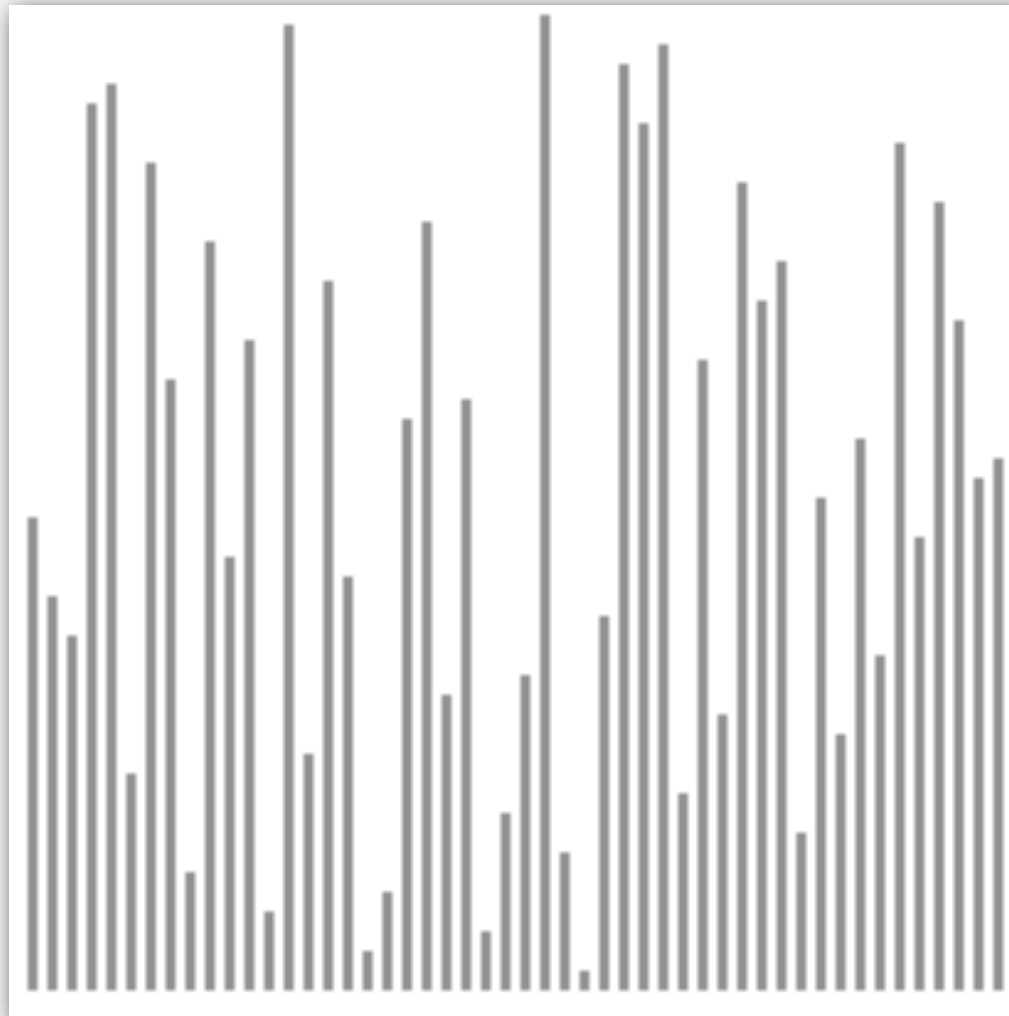magic increment sequence

insertion sort

move to next increment

# Visual trace of shellsort



input

40-sorted

13-sorted

4-sorted

result

# Shellsort animation

http://www.sorting-algorithms.com/shell-sort

▲ algorithm position

■ h-sorted

■ current subsequence

■ other elements

45

# Shellsort animation

http://www.sorting-algorithms.com/shell-sort

▲  algorithm position

█  h-sorted

█  current subsequence

█  other elements

46

## Shellsort: analysis

Proposition. The worst-case number of compares used by shellsort with the 3x+1 increments is $O(N^{3/2})$.

Property. The number of compares used by shellsort with the 3x+1 increments is at most by a small multiple of N times the # of increments used.

| N | compares | $N^{1.289}$ | 2.5 N lg N |
|---|---|---|---|
| 5,000 | 93 | 58 | 106 |
| 10,000 | 209 | 143 | 230 |
| 20,000 | 467 | 349 | 495 |
| 40,000 | 1022 | 855 | 1059 |
| 80,000 | 2266 | 2089 | 2257 |

measured in thousands

Remark. Accurate model has not yet been discovered (!)

Why are we interested in shellsort?

Example of simple idea leading to substantial performance gains.

Useful in practice.

- Fast unless array size is huge.
- Tiny, fixed footprint for code (used in embedded systems).
- Hardware sort prototype.

Simple algorithm, nontrivial performance, interesting questions.

- Asymptotic growth rate?
- Best sequence of increments?  ⟵  open problem:  find a better increment sequence
- Average case performance?

Lesson.  Some good algorithms are still waiting discovery.

# 2.2 Mergesort



▸ mergesort

▸ bottom-up mergesort

▸ sorting complexity

▸ comparators

## Two classic sorting algorithms

Critical components in the world's computational infrastructure.

- Full scientific understanding of their properties has enabled us to develop them into practical system sorts.
- Quicksort honored as one of top 10 algorithms of $20^{th}$ century in science and engineering.

Mergesort.          ⟵ today

- Java sort for objects.
- Perl, Python stable sort.

Quicksort.          ⟵ next lecture

- Java sort for primitive types.
- C qsort, Unix, g++, Visual C++, Python.

# Mergesort

Basic plan.

- Divide array into two halves.
- Recursively sort each half.
- Merge two halves.



| | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| input | M | E | R | G | E | S | O | R | T | E | X | A | M | P | L | E |
| sort left half | E | E | G | M | O | R | R | S | T | E | X | A | M | P | L | E |
| sort right half | E | E | G | M | O | R | R | S | A | E | E | L | M | P | T | X |
| merge results | A | E | E | E | E | G | L | M | M | O | P | R | R | S | T | X |

**Mergesort overview**



First Draft of a Report on the EDVAC

John von Neumann

# Merging

Q. How to combine two sorted subarrays into a sorted whole.

A. Use an auxiliary array.



**Abstract in-place merge trace**

## Merging:  Java implementation

```java
private static void merge(Comparable[] a, int lo, int mid, int hi)
{
   assert isSorted(a, lo, mid);      // precondition: a[lo..mid]   sorted
   assert isSorted(a, mid+1, hi);    // precondition: a[mid+1..hi] sorted

   for (int k = lo; k <= hi; k++)
      aux[k] = a[k];
```

copy

```java
   int i = lo, j = mid+1;
   for (int k = lo; k <= hi; k++)
   {
      if        (i > mid)                a[k] = aux[j++];
      else if (j > hi)                   a[k] = aux[i++];
      else if (less(aux[j], aux[i]))  a[k] = aux[j++];
      else                              a[k] = aux[i++];
   }
```

merge

```java
   assert isSorted(a, lo, hi);      // postcondition: a[lo..hi] sorted
}
```

| | lo | | | i | mid | | | j | | hi |
|--------|---|---|---|---|---|---|---|---|---|---|
| aux[]  | A | G | L | O | R | H | I | M | S | T |

| | | | | | | k | | | |
|------|---|---|---|---|---|---|---|---|---|---|
| a[]  | A | G | H | I | L | M | | | | |

**Assertion.** Statement to test assumptions about your program.

- Helps detect logic bugs.
- Documents code.

**Java assert statement.** Throws an exception unless boolean condition is ture.

```
assert isSorted(a, lo, hi);
```

**Can enable or disable at runtime.** ⇒ No cost in production code.

```
java -ea MyProgram    // enable assertions
java -da MyProgram    // disable assertions (default)
```

**Best practices.** Use to check internal invariants. Assume assertions will be disabled in production code (e.g., don't use for external argument-checking).
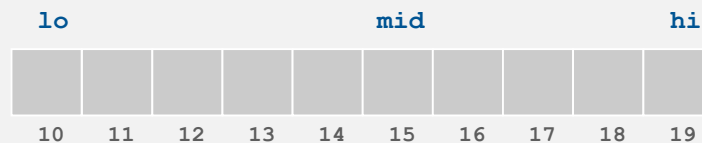
# Mergesort:  Java implementation

```java
public class Merge
{
    private static Comparable[] aux;

    private static void merge(Comparable[] a, int lo, int mid, int hi)
    {  /* as before */  }

    private static void sort(Comparable[] a, int lo, int hi)
    {
        if (hi <= lo) return;
        int mid = lo + (hi - lo) / 2;
        sort(a, lo, mid);
        sort(a, mid+1, hi);
        merge(a, lo, m, hi);
    }

    public static void sort(Comparable[] a)
    {
        aux = new Comparable[a.length];
        sort(a, 0, a.length - 1);
    }
}
```
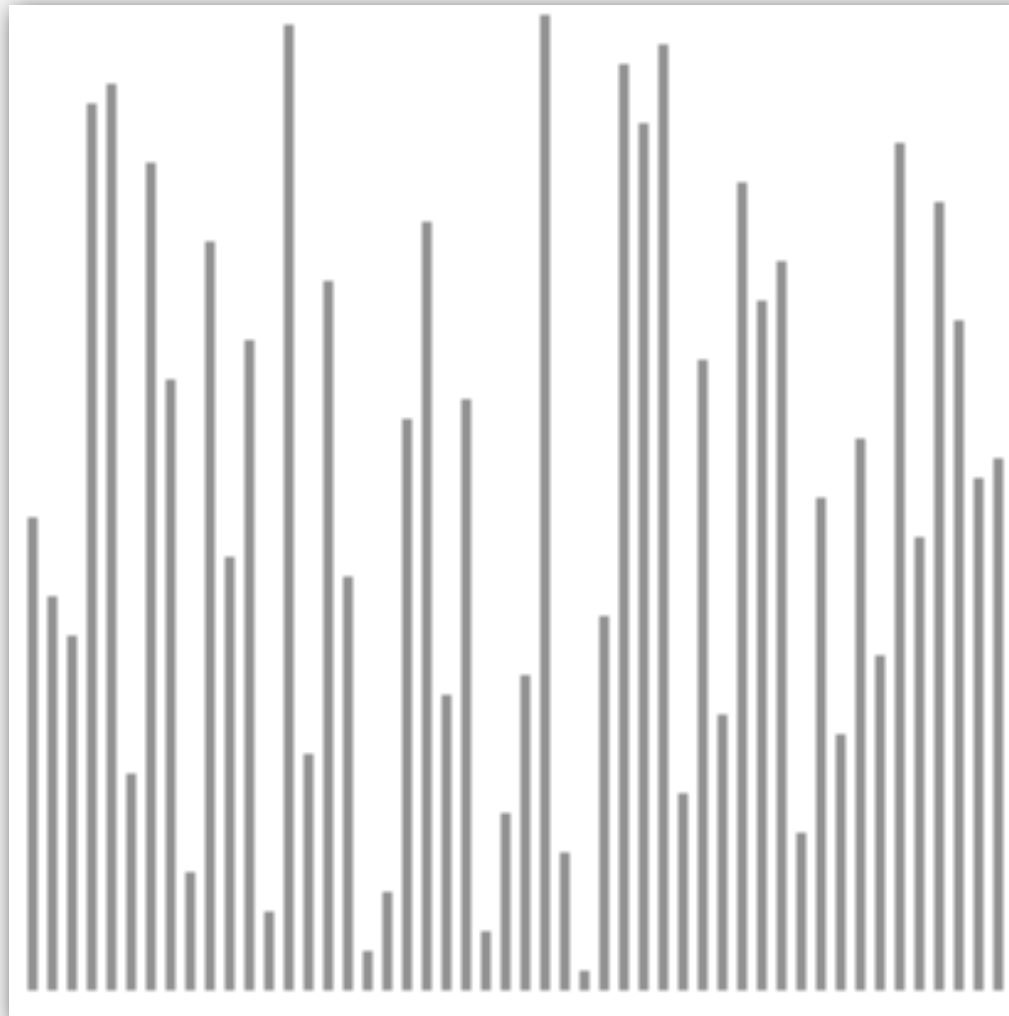
| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| **lo** | | | | | **mid** | | | **hi** | |
| 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |

# Mergesort trace

| | | | a[] | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | lo | | hi | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
| | | | | M | E | R | G | E | S | O | R | T | E | X | A | M | P | L | E |
| merge(a, | 0, | 0, | 1) | E | M | R | G | E | S | O | R | T | E | X | A | M | P | L | E |
| merge(a, | 2, | 2, | 3) | E | M | G | R | E | S | O | R | T | E | X | A | M | P | L | E |
| merge(a, | 0, | 1, | 3) | E | G | M | R | E | S | O | R | T | E | X | A | M | P | L | E |
| merge(a, | 4, | 4, | 5) | E | G | M | R | E | S | O | R | T | E | X | A | M | P | L | E |
| merge(a, | 6, | 6, | 7) | E | G | M | R | E | S | O | R | T | E | X | A | M | P | L | E |
| merge(a, | 4, | 5, | 7) | E | G | M | R | E | O | R | S | T | E | X | A | M | P | L | E |
| merge(a, | 0, | 3, | 7) | E | E | G | M | O | R | R | S | T | E | X | A | M | P | L | E |
| merge(a, | 8, | 8, | 9) | E | E | G | M | O | R | R | S | E | T | X | A | M | P | L | E |
| merge(a, | 10, | 10, | 11) | E | E | G | M | O | R | R | S | E | T | A | X | M | P | L | E |
| merge(a, | 8, | 9, | 11) | E | E | G | M | O | R | R | S | A | E | T | X | M | P | L | E |
| merge(a, | 12, | 12, | 13) | E | E | G | M | O | R | R | S | A | E | T | X | M | P | L | E |
| merge(a, | 14, | 14, | 15) | E | E | G | M | O | R | R | S | A | E | T | X | M | P | E | L |
| merge(a, | 12, | 13, | 15) | E | E | G | M | O | R | R | S | A | E | T | X | E | L | M | P |
| merge(a, | 8, | 11, | 15) | E | E | G | M | O | R | R | S | A | E | E | L | M | P | T | X |
| merge(a, | 0, | 7, | 15) | A | E | E | E | E | G | L | M | M | O | P | R | R | S | T | X |

**Trace of merge results for top-down mergesort**

result after recursive call

9

# Mergesort animation

50 random elements



http://www.sorting-algorithms.com/merge-sort

▲ algorithm position

in order

current subarray

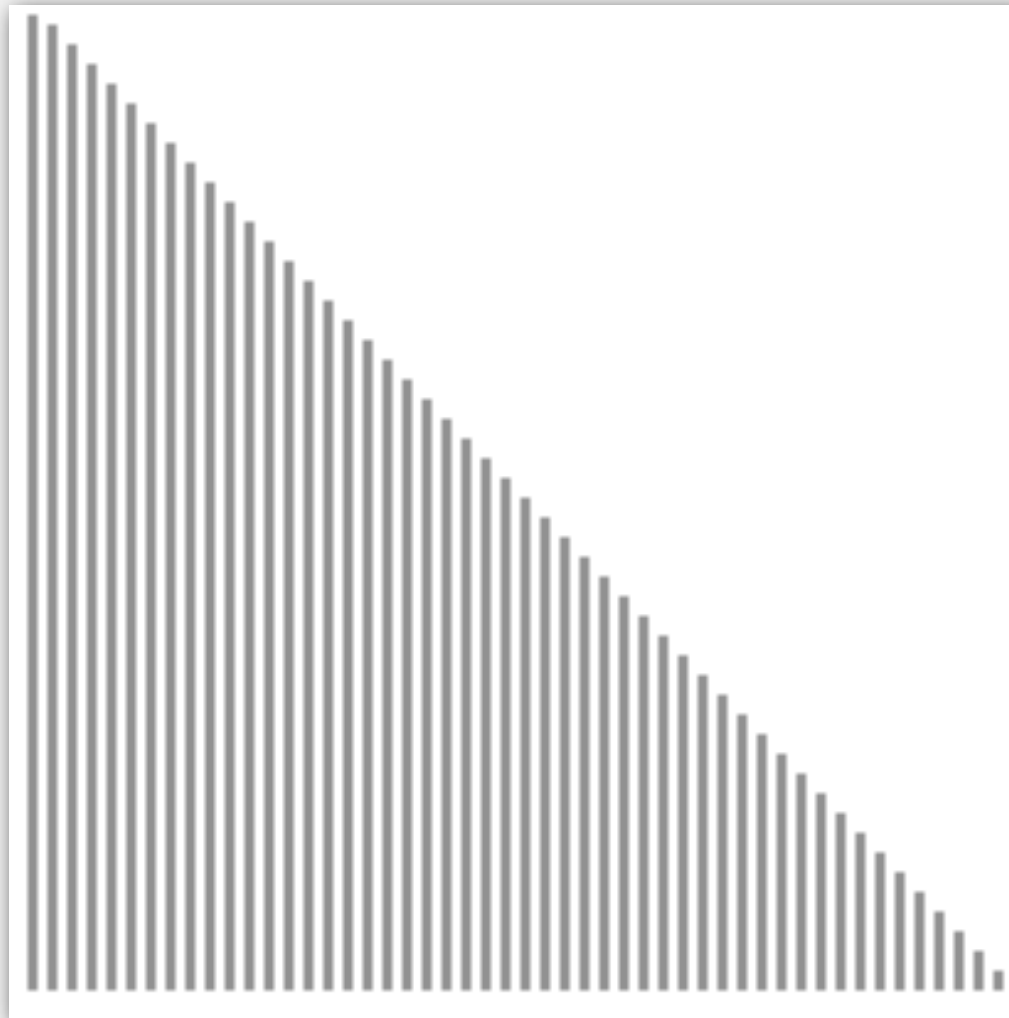not in order

# Mergesort animation

50 reverse-sorted elements



▲ algorithm position

■ in order

■ current subarray

■ not in order

http://www.sorting-algorithms.com/merge-sort

# Mergesort:  empirical analysis

**Running time estimates:**

- Home pc executes $10^8$ comparisons/second.
- Supercomputer executes $10^{12}$ comparisons/second.

| computer | insertion sort ($N^2$) | | | mergesort ($N \log N$) | | |
|---|---|---|---|---|---|---|
| | thousand | million | billion | thousand | million | billion |
| home | instant | 2.8 hours | 317 years | instant | 1 second | 18 min |
| super | instant | 1 second | 1 week | instant | instant | instant |

**Bottom line.**  Good algorithms are better than supercomputers.

## Mergesort: mathematical analysis

**Proposition.** Mergesort uses $\sim 2\,N \lg N$ data moves to sort any array of size $N$.

**Def.** $D(N)$ = number of data moves to mergesort an array of size $N$.

$$= \; D(N/2) \; + \; D(N/2) \; + \; 2\,N$$

left half      right half      merge

**Mergesort recurrence.** $D(N) = 2\,D(N/2) + 2\,N$ for $N > 1$, with $T(1) = 0$.

- Not quite right for odd $N$.
- Similar recurrence holds for many divide-and-conquer algorithms.

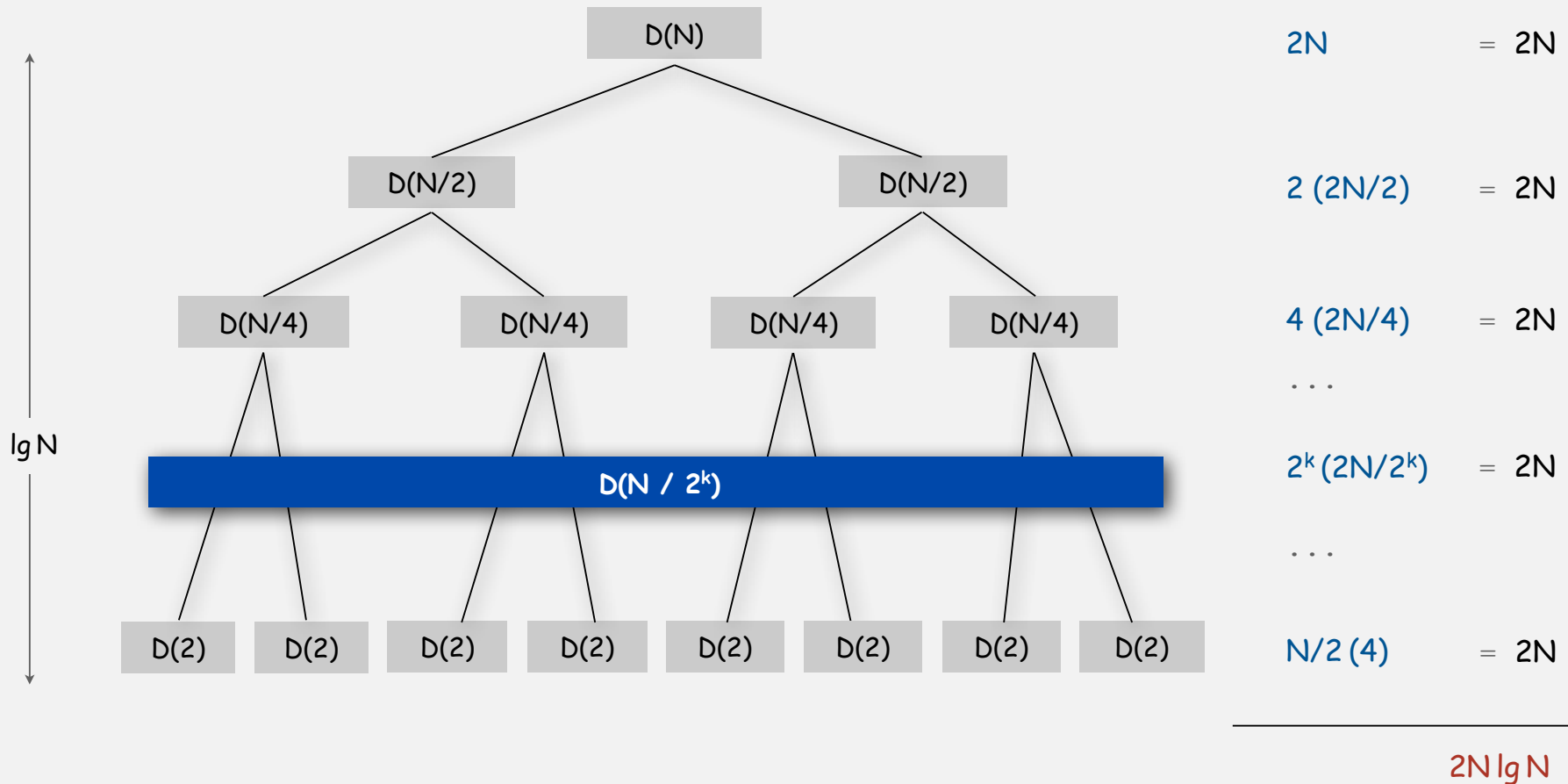**Solution.** $D(N) \sim 2\,N \lg N$.

- For simplicity, we'll prove when $N$ is a power of 2.
- True for all $N$. [see COS 340]

# Mergesort recurrence: proof 1

Mergesort recurrence. $D(N) = 2 D(N/2) + 2 N$ for $N > 1$, with $D(1) = 0$.

Proposition. If $N$ is a power of 2, then $D(N) = 2 N \lg N$.

Pf.



| | |
|---|---|
| 2N | = 2N |
| 2 (2N/2) | = 2N |
| 4 (2N/4) | = 2N |
| ... | |
| $2^k (2N/2^k)$ | = 2N |
| ... | |
| N/2 (4) | = 2N |

$2N \lg N$

## Mergesort recurrence: proof 2

Mergesort recurrence. $D(N) = 2 \, D(N / 2) + 2 \, N$ for $N > 1$, with $D(1) = 0$.

Proposition. If $N$ is a power of 2, then $D(N) = 2 \, N \lg N$.

Pf.

| | |
|---|---|
| D(N)  = 2 D(N/2) + 2N | given |
| D(N) / N = 2 D(N/2) / N + 2 | divide both sides by N |
| = D(N/2) / (N/2) + 2 | algebra |
| = D(N/4) / (N/4) + 2 + 2 | apply to first term |
| = D(N/8) / (N/8) + 2 + 2 + 2 | apply to first term again |
| . . . | |
| = D(N/N) / (N/N) + 2 + 2 + ... + 2 | stop applying, T(1) = 0 |
| = 2 lg N | |

**Mergesort recurrence.**  $D(N) = 2\,D(N/2) + 2\,N$  for $N > 1$, with $D(1) = 0$.

**Proposition.**  If $N$ is a power of 2, then $D(N) = 2\,N \lg N$.

**Pf.**  [by induction on N]

• **Base case:**  $N = 1$.

• **Inductive hypothesis:**  $D(N) = 2N \lg N$.

• **Goal:**  show that $D(2N) = 2(2N)\lg(2N)$.

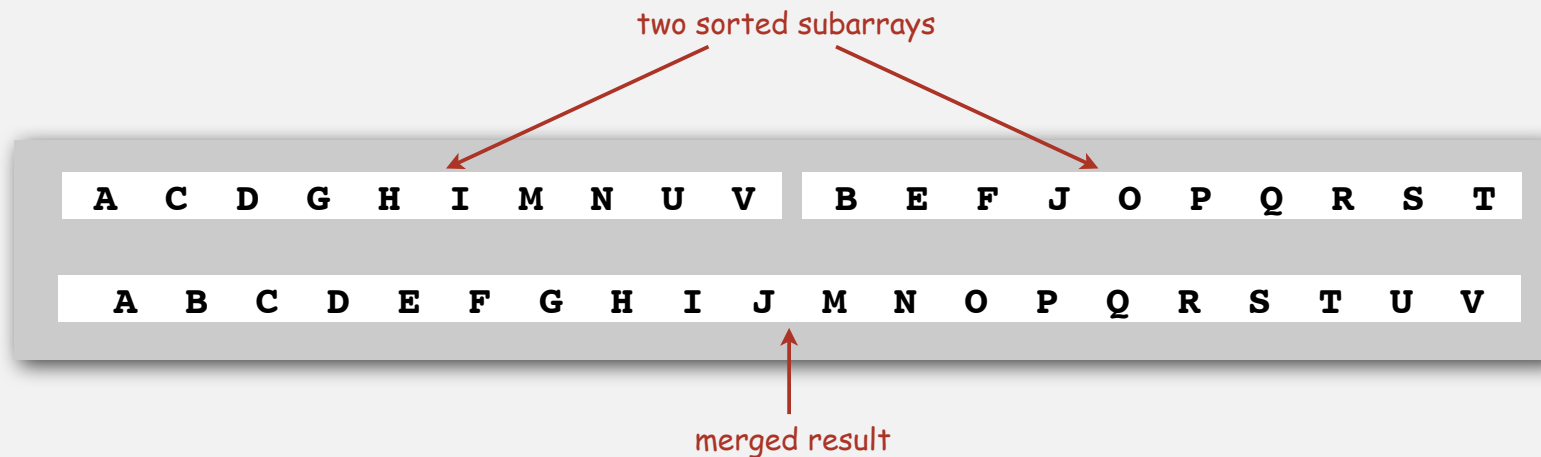| | |
|---|---|
| $D(2N) \;=\; 2\,D(N) \;+\; 4N$ | given |
| $=\; 4\,N \lg N + 4\,N$ | inductive hypothesis |
| $=\; 4\,N\,(\lg(2N) - 1) + 4N$ | algebra |
| $=\; 4\,N \lg(2N)$ | QED |

Mergesort:  number of compares

Proposition.  Mergesort uses between $\frac{1}{2} N \lg N$ and $N \lg N$ compares to sort any array of size $N$.

Pf.  The number of compares for the last merge is between $\frac{1}{2} N \lg N$ and $N$.

## Mergesort analysis: memory

Proposition G. Mergesort uses extra space proportional to N.

Pf. The array `aux[]` needs to be of size N for the last merge.

two sorted subarrays

| A | C | D | G | H | I | M | N | U | V | | B | E | F | J | O | P | Q | R | S | T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| A | B | C | D | E | F | G | H | I | J | M | N | O | P | Q | R | S | T | U | V |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

merged result

Def. A sorting algorithm is in-place if it uses O(log N) extra memory.

Ex. Insertion sort, selection sort, shellsort.

Challenge for the bored. In-place merge. [Kronrud, 1969]

## Mergesort: practical improvements

### Use insertion sort for small subarrays.

- Mergesort has too much overhead for tiny subarrays.
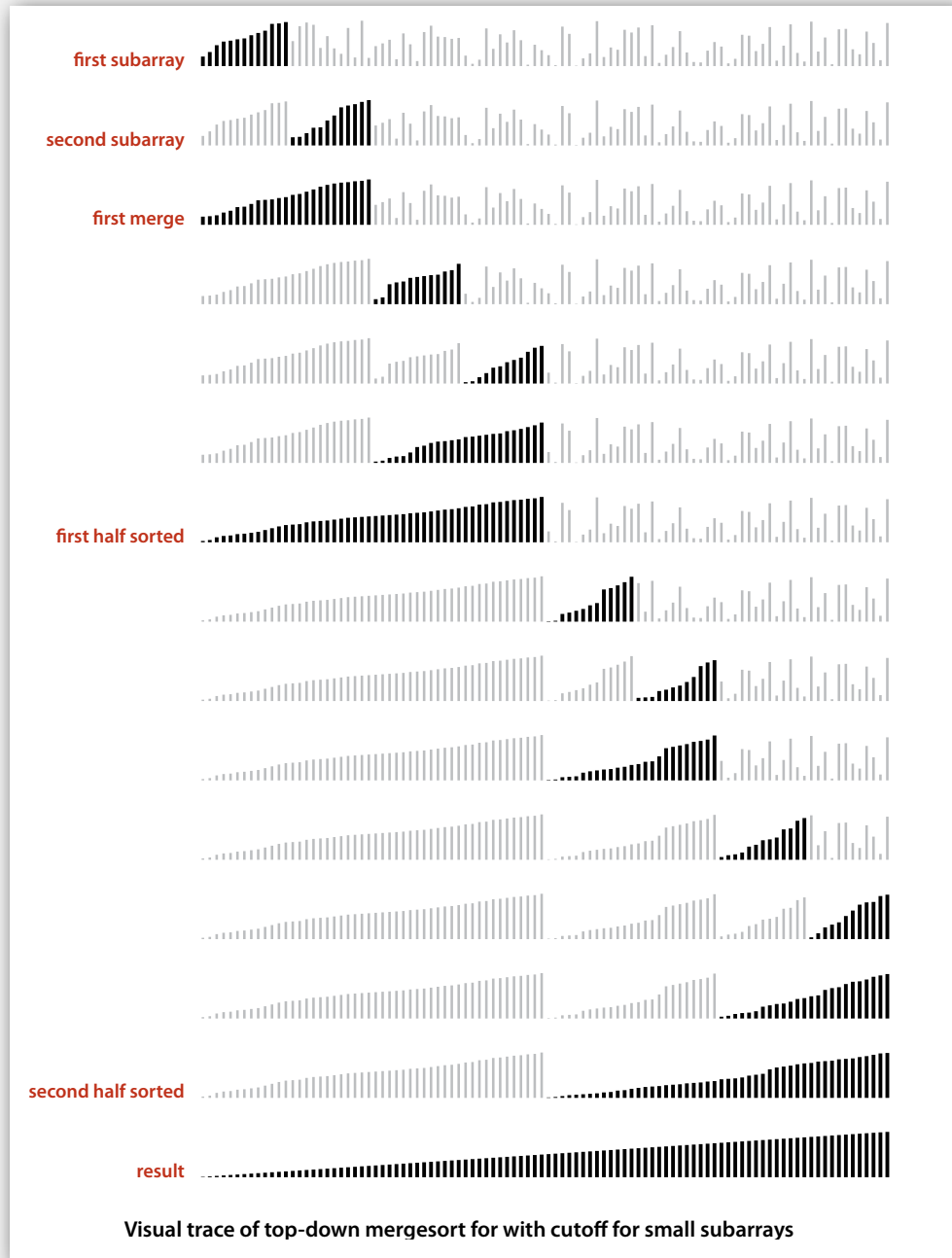- Cutoff to insertion sort for ≈ 7 elements.

### Stop if already sorted.

- Is biggest element in first half ≤ smallest element in second half?
- Helps for partially-ordered arrays.

| A | B | C | D | E | F | G | H | I | J | M | N | O | P | Q | R | S | T | U | V |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| A | B | C | D | E | F | G | H | I | J | M | N | O | P | Q | R | S | T | U | V |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

### Eliminate the copy to the auxiliary array. Save time (but not space)
by switching the role of the input and auxiliary array in each recursive call.

### Ex. See `MergeX.java` or `Arrays.sort()`.

# Mergesort visualization



first subarray

second subarray

first merge

first half sorted

second half sorted

result

**Visual trace of top-down mergesort for with cutoff for small subarrays**

# Bottom-up mergesort

## Basic plan.

- Pass through array, merging subarrays of size 1.
- Repeat for subarrays of size 2, 4, 8, 16, ....

```
                                               a[i]
                          0   1   2   3   4   5   6   7   8   9  10  11  12  13  14  15
                          M   E   R   G   E   S   O   R   T   E   X   A   M   P   L   E
   sz = 2
   merge(a,   0,   0,   1)  E   M   R   G   E   S   O   R   T   E   X   A   M   P   L   E
   merge(a,   2,   2,   3)  E   M   G   R   E   S   O   R   T   E   X   A   M   P   L   E
   merge(a,   4,   4,   5)  E   M   G   R   E   S   O   R   T   E   X   A   M   P   L   E
   merge(a,   6,   6,   7)  E   M   G   R   E   S   O   R   T   E   X   A   M   P   L   E
   merge(a,   8,   8,   9)  E   M   G   R   E   S   O   R   E   T   X   A   M   P   L   E
   merge(a,  10,  10,  11)  E   M   G   R   E   S   O   R   E   T   A   X   M   P   L   E
   merge(a,  12,  12,  13)  E   M   G   R   E   S   O   R   E   T   A   X   M   P   L   E
   merge(a,  14,  14,  15)  E   M   G   R   E   S   O   R   E   T   A   X   M   P   E   L
  sz = 4
   merge(a,   0,   1,   3)  E   G   M   R   E   S   O   R   E   T   A   X   M   P   E   L
   merge(a,   4,   5,   7)  E   G   M   R   E   O   R   S   E   T   A   X   M   P   E   L
   merge(a,   8,   9,  11)  E   G   M   R   E   O   R   S   A   E   T   X   M   P   E   L
   merge(a,  12,  13,  15)  E   G   M   R   E   O   R   S   A   E   T   X   E   L   M   P
  sz = 8
   merge(a,   0,   3,   7)  E   E   G   M   O   R   R   S   A   E   T   X   E   L   M   P
   merge(a,   8,  11,  15)  E   E   G   M   O   R   R   S   A   E   E   L   M   P   T   X
 sz = 16
   merge(a,   0,   7,  15)  A   E   E   E   E   G   L   M   M   O   P   R   R   S   T   X
```

**Trace of merge results for bottom-up mergesort**

## Bottom line. No recursion needed!

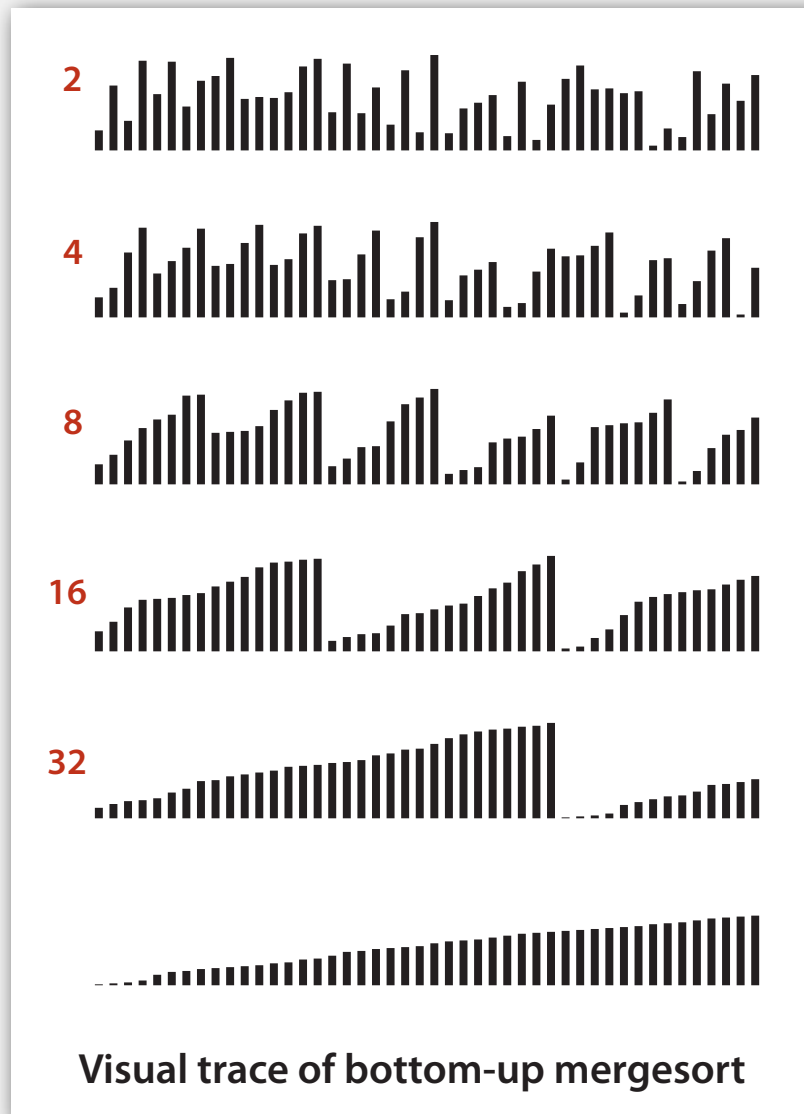## Bottom-up mergesort: Java implementation

```java
public class MergeBU
{
   private static Comparable[] aux;

   private static void merge(Comparable[] a, int lo, int mid, int hi)
   {  /* as before */  }

   public static void sort(Comparable[] a)
   {
      int N = a.length;
      aux = new Comparable[N];
      for (int sz = 1; sz < N; sz = sz+sz)
         for (int lo = 0; lo < N-sz; lo += sz+sz)
            merge(a, lo, lo+sz-1, Math.min(lo+sz+sz-1, N-1));
   }
}
```

Bottom line.  Concise industrial-strength code, if you have the space.

# Bottom-up mergesort:  visual trace



**Visual trace of bottom-up mergesort**

## Complexity of sorting

Computational complexity.  Framework to study efficiency of algorithms for solving a particular problem X.

Machine model.  Focus on fundamental operations.

Upper bound.  Cost guarantee provided by some algorithm for X.

Lower bound.  Proven limit on cost guarantee of all algorithms for X.

Optimal algorithm.  Algorithm with best cost guarantee for X.

lower bound ~ upper bound

access information only through compares

Example:  sorting.

- Machine model = # compares.
- Upper bound = ~ N lg N from mergesort.
- Lower bound = ~ N lg N ?
- Optimal algorithm = mergesort ?

# Decision tree (for 3 distinct elements)

## Compare-based lower bound for sorting

Proposition. Any compare-based sorting algorithm must use at least $\lg N! \sim N \lg N$ compares in the worst-case.

Pf.

- Assume input consists of $N$ distinct values $a_1$ through $a_N$.
- Worst case dictated by height $h$ of decision tree.
- Binary tree of height $h$ has at most $2^h$ leaves.
- $N!$ different orderings $\Rightarrow$ at least $N!$ leaves.



at least N! leaves

no more than $2^h$ leaves

$h$

## Compare-based lower bound for sorting

Proposition. Any compare-based sorting algorithm must use at least $\lg N! \sim N \lg N$ compares in the worst-case.

Pf.

- Assume input consists of $N$ distinct values $a_1$ through $a_N$.
- Worst case dictated by height $h$ of decision tree.
- Binary tree of height $h$ has at most $2^h$ leaves.
- $N!$ different orderings $\Rightarrow$ at least $N!$ leaves.

$$2^h \geq \# \text{ leaves} \geq N!$$
$$\Rightarrow h \geq \lg N! \sim N \lg N$$

Stirling's formula

## Complexity of sorting

Machine model. Focus on fundamental operations.

Upper bound. Cost guarantee provided by some algorithm for X.

Lower bound. Proven limit on cost guarantee of all algorithms for X.

Optimal algorithm. Algorithm with best cost guarantee for X.

Example: sorting.

- Machine model = # compares.
- Upper bound = ~ N lg N from mergesort.
- Lower bound = ~ N lg N.
- Optimal algorithm = mergesort.

First goal of algorithm design: optimal algorithms.

## Complexity results in context

Other operations?  Mergesort optimality is only about number of compares.

Space?
- Mergesort is not optimal with respect to space usage.
- Insertion sort, selection sort, and shellsort are space-optimal.

Challenge.  Find an algorithm that is both time- and space-optimal.

Lessons.  Use theory as a guide.

Ex.  Don't try to design sorting algorithm that uses $\frac{1}{2} N \lg N$ compares.

## Complexity results in context (continued)

Lower bound may not hold if the algorithm has information about:
- The initial order of the input.
- The distribution of key values.
- The representation of the keys.

**Partially-ordered arrays.** Depending on the initial order of the input,

we may not need N lg N compares.

*insertion sort requires only N-1
compares on an already sorted array*

**Duplicate keys.** Depending on the input distribution of duplicates,

we may not need N lg N compares.

*stay tuned for 3-way quicksort*

**Digital properties of keys.** We can use digit/character compares instead of

key compares for numbers and strings.

*stay tuned for radix sorts*

33

# Sort by artist name



| | Name | Artist | Time | Album |
|---|---|---|---|---|
| 12 | ☑ Let It Be | The Beatles | 4:03 | Let It Be |
| 13 | ☑ Take My Breath Away | BERLIN | 4:13 | Top Gun – Soundtrack |
| 14 | ☑ Circle Of Friends | Better Than Ezra | 3:27 | Empire Records |
| 15 | ☑ Dancing With Myself | Billy Idol | 4:43 | Don't Stop |
| 16 | ☑ Rebel Yell | Billy Idol | 4:49 | Rebel Yell |
| 17 | ☑ Piano Man | Billy Joel | 5:36 | Greatest Hits Vol. 1 |
| 18 | ☑ Pressure | Billy Joel | 3:16 | Greatest Hits, Vol. II (1978 – 1985) (Disc 2) |
| 19 | ☑ The Longest Time | Billy Joel | 3:36 | Greatest Hits, Vol. II (1978 – 1985) (Disc 2) |
| 20 | ☑ Atomic | Blondie | 3:50 | Atomic: The Very Best Of Blondie |
| 21 | ☑ Sunday Girl | Blondie | 3:15 | Atomic: The Very Best Of Blondie |
| 22 | ☑ Call Me | Blondie | 3:33 | Atomic: The Very Best Of Blondie |
| 23 | ☑ Dreaming | Blondie | 3:06 | Atomic: The Very Best Of Blondie |
| 24 | ☑ Hurricane | Bob Dylan | 8:32 | Desire |
| 25 | ☑ The Times They Are A–Changin' | Bob Dylan | 3:17 | Greatest Hits |
| 26 | ☑ Livin' On A Prayer | Bon Jovi | 4:11 | Cross Road |
| 27 | ☑ Beds Of Roses | Bon Jovi | 6:35 | Cross Road |
| 28 | ☑ Runaway | Bon Jovi | 3:53 | Cross Road |
| 29 | ☑ Rasputin (Extended Mix) | Boney M | 5:50 | Greatest Hits |
| 30 | ☑ Have You Ever Seen The Rain | Bonnie Tyler | 4:10 | Faster Than The Speed Of Night |
| 31 | ☑ Total Eclipse Of The Heart | Bonnie Tyler | 7:02 | Faster Than The Speed Of Night |
| 32 | ☑ Straight From The Heart | Bonnie Tyler | 3:41 | Faster Than The Speed Of Night |
| 33 | ☑ Holding Out For A Hero | Bonny Tyler | 5:49 | Meat Loaf And Friends |
| 34 | ☑ Dancing In The Dark  ⊙ | Bruce Springsteen ⊙ | 4:05 | Born In The U.S.A. |
| 35 | ☑ Thunder Road | Bruce Springsteen | 4:51 | Born To Run |
| 36 | ☑ Born To Run | Bruce Springsteen | 4:30 | Born To Run |
| 37 | ☑ Jungleland | Bruce Springsteen | 9:34 | Born To Run |
| 38 | ☑ Turn! Turn! Turn! (To Everythin… | The Byrds | 3:57 | Forrest Gump The Soundtrack (Disc 2) |

34

# Sort by song name

## Natural order

Comparable interface: sort uses type's natural order.

```java
public class Date implements Comparable<Date>
{
   private final int month, day, year;

   public Date(int m, int d, int y)
   {
      month = m;
      day   = d;
      year  = y;
   }

   …
   public int compareTo(Date that)
   {
      if (this.year  < that.year ) return -1;
      if (this.year  > that.year ) return +1;
      if (this.month < that.month) return -1;
      if (this.month > that.month) return +1;
      if (this.day   < that.day  ) return -1;
      if (this.day   > that.day  ) return +1;
      return 0;
   }
}
```

natural order

## Generalized compare

Comparable interface: sort uses type's natural order.

Problem 1. May want to use a non-natural order.

Problem 2. Desired data type may not come with a "natural" order.

Ex. Sort strings by:

- Natural order.        `Now is the time`
- Case insensitive.     `is Now the time`
- Spanish.              `café cafetero cuarto churro nube ñoño`
- British phone book.   `McKinley Mackintosh`

pre-1994 order for digraphs
ch and ll and rr

```
String[] a;
...
Arrays.sort(a);
Arrays.sort(a, String.CASE_INSENSITIVE_ORDER);
Arrays.sort(a, Collator.getInstance(Locale.SPANISH));
```

import `java.text.Collator;`

## Comparators

Solution. Use Java's `Comparator` interface.

```
public interface Comparator<Key>
{
    public int compare(Key v, Key w);
}
```

Remark. The `compare()` method implements a total order like `compareTo()`.

Advantages. Decouples the definition of the data type from the definition of what it means to compare two objects of that type.
- Can add any number of new orders to a data type.
- Can add an order to a library data type with no natural order.

## Comparator example

**Reverse order.** Sort an array of strings in reverse order.

```java
public class ReverseOrder implements Comparator<String>
{
    public int compare(String a, String b)
    {
        return b.compareTo(a);
    }
}
```

comparator implementation

```java
    ...
    Arrays.sort(a, new ReverseOrder());
    ...
```

client

## Sort implementation with comparators

To support comparators in our sort implementations:

- Pass `Comparator` to `sort()` and `less()`.
- Use it in `less()`.

Ex.  Insertion sort.

```
public static void sort(Object[] a, Comparator comparator)
{
   int N = a.length;
   for (int i = 0; i < N; i++)
      for (int j = i; j > 0 && less(comparator, a[j], a[j-1]); j--)
         exch(a, j, j-1);
}

private static boolean less(Comparator c, Object v, Object w)
{   return c.compare(v, w) < 0;    }

private static void exch(Object[] a, int i, int j)
{   Object swap = a[i]; a[i] = a[j]; a[j] = swap;   }
```

## Generalized compare

Comparators enable multiple sorts of a single array (by different keys).

Ex. Sort students by name or by section.

```
Arrays.sort(students, Student.BY_NAME);
Arrays.sort(students, Student.BY_SECT);
```

sort by name

| Andrews | 3 | A | 664-480-0023 | 097 Little |
|---------|---|---|--------------|------------|
| Battle | 4 | C | 874-088-1212 | 121 Whitman |
| Chen | 2 | A | 991-878-4944 | 308 Blair |
| Fox | 1 | A | 884-232-5341 | 11 Dickinson |
| Furia | 3 | A | 766-093-9873 | 101 Brown |
| Gazsi | 4 | B | 665-303-0266 | 22 Brown |
| Kanaga | 3 | B | 898-122-9643 | 22 Brown |
| Rohde | 3 | A | 232-343-5555 | 343 Forbes |

sort by section

| Fox | 1 | A | 884-232-5341 | 11 Dickinson |
|---------|---|---|--------------|------------|
| Chen | 2 | A | 991-878-4944 | 308 Blair |
| Andrews | 3 | A | 664-480-0023 | 097 Little |
| Furia | 3 | A | 766-093-9873 | 101 Brown |
| Kanaga | 3 | B | 898-122-9643 | 22 Brown |
| Rohde | 3 | A | 232-343-5555 | 343 Forbes |
| Battle | 4 | C | 874-088-1212 | 121 Whitman |
| Gazsi | 4 | B | 665-303-0266 | 22 Brown |

## Generalized compare

Ex. Enable sorting students by name or by section.

```
public class Student
{
   public static final Comparator<Student> BY_NAME = new ByName();
   public static final Comparator<Student> BY_SECT = new BySect();

   private final String name;
   private final int section;
   ...
   private static class ByName implements Comparator<Student>
   {
      public int compare(Student a, Student b)
      {   return a.name.compareTo(b.name);   }
   }

   private static class BySect implements Comparator<Student>
   {
      public int compare(Student a, Student b)
      {   return a.section - b.section;   }
   }
}
```

only use this trick if no danger of overflow

## Generalized compare problem

A typical application.  First, sort by name; then sort by section.

`Arrays.sort(students, Student.BY_NAME);`

| Andrews | 3 | A | 664-480-0023 | 097 Little |
|---------|---|---|--------------|------------|
| Battle  | 4 | C | 874-088-1212 | 121 Whitman |
| Chen    | 2 | A | 991-878-4944 | 308 Blair |
| Fox     | 1 | A | 884-232-5341 | 11 Dickinson |
| Furia   | 3 | A | 766-093-9873 | 101 Brown |
| Gazsi   | 4 | B | 665-303-0266 | 22 Brown |
| Kanaga  | 3 | B | 898-122-9643 | 22 Brown |
| Rohde   | 3 | A | 232-343-5555 | 343 Forbes |

`Arrays.sort(students, Student.BY_SECT);`

| Fox     | 1 | A | 884-232-5341 | 11 Dickinson |
|---------|---|---|--------------|------------|
| Chen    | 2 | A | 991-878-4944 | 308 Blair |
| Kanaga  | 3 | B | 898-122-9643 | 22 Brown |
| Andrews | 3 | A | 664-480-0023 | 097 Little |
| Furia   | 3 | A | 766-093-9873 | 101 Brown |
| Rohde   | 3 | A | 232-343-5555 | 343 Forbes |
| Battle  | 4 | C | 874-088-1212 | 121 Whitman |
| Gazsi   | 4 | B | 665-303-0266 | 22 Brown |

@#%&@!!.  Students in section 3 no longer in order by name.

A stable sort preserves the relative order of records with equal keys.

# Sorting challenge 5

Q. Which sorts are stable?

Insertion sort?  Selection sort?  Shellsort?  Mergesort?

| sorted by time | sorted by location (not stable) | sorted by location (stable) |
|---|---|---|
| Chicago   09:00:00 | Chicago 09:25:52 | Chicago 09:00:00 |
| Phoenix   09:00:03 | Chicago 09:03:13 | Chicago 09:00:59 |
| Houston   09:00:13 | Chicago 09:21:05 | Chicago 09:03:13 |
| Chicago   09:00:59 | Chicago 09:19:46 | Chicago 09:19:32 |
| Houston   09:01:10 | Chicago 09:19:32 | Chicago 09:19:46 |
| Chicago   09:03:13 | Chicago 09:00:00 | Chicago 09:21:05 |
| Seattle   09:10:11 | Chicago 09:35:21 | Chicago 09:25:52 |
| Seattle   09:10:25 | Chicago 09:00:59 | Chicago 09:35:21 |
| Phoenix   09:14:25 | Houston 09:01:10 | Houston 09:00:13 |
| Chicago   09:19:32 | Houston 09:00:13 | Houston 09:01:10 |
| Chicago   09:19:46 | Phoenix 09:37:44 | Phoenix 09:00:03 |
| Chicago   09:21:05 | Phoenix 09:00:03 | Phoenix 09:14:25 |
| Seattle   09:22:43 | Phoenix 09:14:25 | Phoenix 09:37:44 |
| Seattle   09:22:54 | Seattle 09:10:25 | Seattle 09:10:11 |
| Chicago   09:25:52 | Seattle 09:36:14 | Seattle 09:10:25 |
| Chicago   09:35:21 | Seattle 09:22:43 | Seattle 09:22:43 |
| Seattle   09:36:14 | Seattle 09:10:11 | Seattle 09:22:54 |
| Phoenix   09:37:44 | Seattle 09:22:54 | Seattle 09:36:14 |

*no longer sorted by time*

*still sorted by time*

**Stability when sorting on a second key**

Q. Is insertion sort stable?

```
public class Insertion
{
   public static void sort(Comparable[] a)
   {
      int N = a.length;
      for (int i = 0; i < N; i++)
         for (int j = i; j > 0 && less(a[j], a[j-1]); j--)
            exch(a, j, j-1);
   }
}
```

| i | j | 0 | 1 | 2 | 3 | 4 |
|---|---|-----|-----|-----|-----|-----|
| 0 | 0 | $B_1$ | $A_1$ | $A_2$ | $A_3$ | $B_2$ |
| 1 | 0 | $A_1$ | $B_1$ | $A_2$ | $A_3$ | $B_2$ |
| 2 | 1 | $A_1$ | $A_2$ | $B_1$ | $A_3$ | $B_2$ |
| 3 | 2 | $A_1$ | $A_2$ | $A_3$ | $B_1$ | $B_2$ |
| 4 | 4 | $A_1$ | $A_2$ | $A_3$ | $B_1$ | $B_2$ |
|   |   | $A_1$ | $A_2$ | $A_3$ | $B_1$ | $B_2$ |

A. Yes, equal elements never more past each other.

Q. Is selection sort stable ?

```
public class Selection
{
   public static void sort(Comparable[] a)
   {
      int N = a.length;
      for (int i = 0; i < N; i++)
      {
         int min = i;
         for (int j = i+1; j < N; j++)
            if (less(a[j], a[min]))
               min = j;
         exch(a, i, min);
      }
   }
}
```

| i | min | 0 | 1 | 2 |
|---|-----|---|---|---|
| 0 | 2 | $B_1$ | $B_2$ | A |
| 1 | 1 | A | $B_2$ | $B_1$ |
| 2 | 2 | A | $B_2$ | $B_1$ |
|   |   | A | $B_2$ | $B_1$ |

A. No, long-distance exchange might move left element to the right of some equal element.

Q. Is shellsort stable?

```
public class Shell
{
   public static void sort(Comparable[] a)
   {
      int N = a.length;
      int h = 1;
      while (h < N/3) h = 3*h + 1;
      while (h >= 1)
      {
         for (int i = h; i < N; i++)
         {
            for (int j = i; j > h && less(a[j], a[j-h]); j -= h)
               exch(a, j, j-h);
         }
         h = h/3;
      }
   }
}
```

| h | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
|   | $B_1$ | $B_2$ | $B_3$ | $B_4$ | $A_1$ |
| 4 | $A_1$ | $B_2$ | $B_3$ | $B_4$ | $B_1$ |
| 1 | $A_1$ | $B_2$ | $B_3$ | $B_4$ | $B_1$ |
|   | $A_1$ | $B_2$ | $B_3$ | $B_4$ | $B_1$ |

A. No. Long-distance exchanges.

Q. Is mergesort stable?

```java
public class Merge
{
   private static Comparable[] aux;
   private static void merge(Comparable[] a, int lo, int mid, int hi)
   {  /* as before */  }

   private static void sort(Comparable[] a, int lo, int hi)
   {
      if (hi <= lo) return;
      int mid = lo + (hi - lo) / 2;
      sort(a, lo, mid);
      sort(a, mid+1, hi);
      merge(a, lo, mid, hi);
   }

   public static void sort(Comparable[] a)
   {
      aux = new Comparable[a.length];
      sort(a, 0, a.length - 1);
   }
}
```

# Sorting challenge 5D

**Q.** Is mergesort stable?

|  | lo | m | hi | a[i] |
|---|---|---|---|---|

|  | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | M | E | R | G | E | S | O | R | T | E | X | A | M | P | L | E |
| merge(a, 0, 0, 1) | E | M | R | G | E | S | O | R | T | E | X | A | M | P | L | E |
| merge(a, 2, 2, 3) | E | M | G | R | E | S | O | R | T | E | X | A | M | P | L | E |
| merge(a, 4, 4, 5) | E | G | M | R | E | S | O | R | T | E | X | A | M | P | L | E |
| merge(a, 6, 6, 7) | E | G | M | R | E | S | O | R | T | E | X | A | M | P | L | E |
| merge(a, 8, 8, 9) | E | E | G | M | O | R | R | S | E | T | X | A | M | P | L | E |
| merge(a, 10, 10, 11) | E | E | G | M | O | R | R | S | E | T | A | X | M | P | L | E |
| merge(a, 12, 12, 13) | E | E | G | M | O | R | R | S | A | E | T | X | M | P | L | E |
| merge(a, 14, 14, 15) | E | E | G | M | O | R | R | S | A | E | T | X | M | P | E | L |
| merge(a, 0, 1, 3) | E | G | M | R | E | S | O | R | T | E | X | A | M | P | L | E |
| merge(a, 4, 5, 7) | E | G | M | R | E | O | R | S | T | E | X | A | M | P | L | E |
| merge(a, 8, 9, 11) | E | E | G | M | O | R | R | S | A | E | T | X | M | P | L | E |
| merge(a, 12, 13, 15) | E | E | G | M | O | R | R | S | A | E | T | X | E | L | M | P |
| merge(a, 0, 3, 7) | E | E | G | M | O | R | R | S | T | E | X | A | M | P | L | E |
| merge(a, 8, 11, 15) | E | E | G | M | O | R | R | S | A | E | E | L | M | P | T | X |
| merge(a, 0, 7, 15) | A | E | E | E | E | G | L | M | M | O | P | R | R | S | T | X |

**Trace of merge results for bottom-up mergesort**

**A.** Yes, if merge is stable.

Q. Is merge stable?

```
private static void merge(Comparable[] a, int lo, int mid, int hi)
{
    for (int k = lo; k <= hi; k++)
        aux[k] = a[k];

    int i = lo, j = mid+1;
    for (int k = lo; k <= hi; k++)
    {
        if       (i > mid)              a[k] = aux[j++];
        else if (j > hi)                a[k] = aux[i++];
        else if (less(aux[j], aux[i]))  a[k] = aux[j++];
        else                            a[k] = aux[i++];
    }
}
```
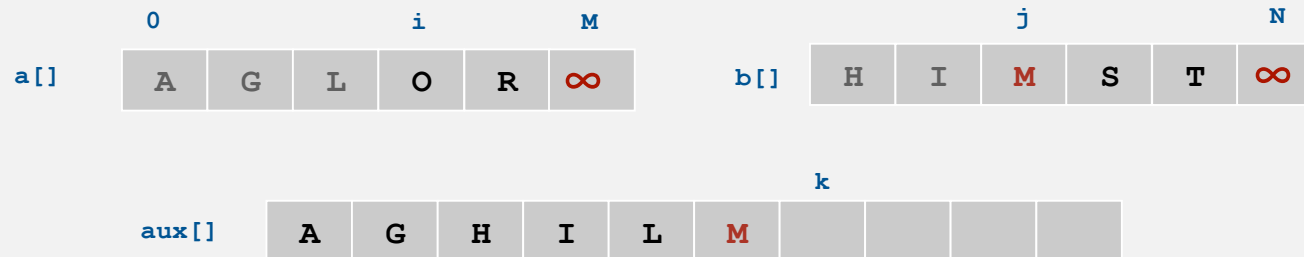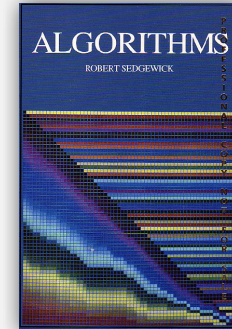
A. Yes, if implemented carefully (take from left subarray if equal).

Sorting challenge 5 (summary)

Q.  Which sorts are stable ?

Yes.  Insertion sort, mergesort.
No.  Selection sort, shellsort.

Note.  Need to carefully check code ("less than" vs "less than or equal").

## Postscript: optimizing mergesort (a short history)

Goal. Remove instructions from the inner loop.

```
private static void merge(Comparable[] a, int lo, int mid, int hi)
{

    for (int k = lo; k <= hi; k++)
        aux[k] = a[k];

    int i = lo, j = mid+1;
    for (int k = lo; k <= hi; k++)
        if        (i > mid)              a[k] = aux[j++];
        else if (j > hi )                a[k] = aux[i++];
        else if (less(aux[j], aux[i]))   a[k] = aux[j++];
        else                             a[k] = aux[i++];

}
```

|  | lo |  |  | i | mid |  |  | j |  | hi |
|---|---|---|---|---|---|---|---|---|---|---|
| aux[] | A | G | L | O | R | H | I | M | S | T |

|  |  |  |  |  | k |  |  |  |  |  |
|---|---|---|---|---|---|---|---|---|---|---|
| a[] | A | G | H | I | L | M |  |  |  |  |

## Postscript: optimizing mergesort (a short history)

**Idea 1 (1960s).** Use sentinels.

```
a[M] := maxint; b[N] := maxint;
for (int i = 0, j = 0, k = 0; k < M+1; k++)
    if (less(aux[j], aux[i])) aux[k] = a[i++];
                              aux[k] = b[j++];
```

Problem 1.  Still need copy.

Problem 2. No good place to put sentinels.

Problem 3. Complicates data-type interface (what is infinity for your type?)

# Postscript: Optimizing mergesort (a short history)

**Idea 2 (1980s).** Reverse copy.

```
private static void merge(Comparable[] a, int lo, int mid, int hi)
{
   for (int i = lo; i <= mid; i++)
      aux[i] = a[i];                                              copy

   for (int j = mid+1; j <= hi; j++)
      aux[j] = a[hi-j+mid+1];                              reverse copy

   int i = lo, j = hi;
   for (int k = lo; k <= hi; k++)
      if (less(aux[j], aux[i])) a[k] = aux[j--];               merge
      else                      a[k] = aux[i++];
}
```

|       | lo |   |   | i | mid | j |   |   | hi |
|-------|----|---|---|---|-----|---|---|---|----|
| aux[] | A  | G | L | O | R   | T | S | M | I  | H |

|      |   |   |   | k |   |   |
|------|---|---|---|---|---|---|
| a[]  | A | G | H | I | L | M |

**Problem.** Copy still in inner loop.

55

# Postscript: Optimizing mergesort (a short history)

Idea 3 (1990s). Eliminate copy with recursive argument switch.

```
int mid = (lo+hi)/2;
mergesortABr(b, a, lo, mid);
mergesortABr(b, a, mid+1, r);
mergeAB(a, lo, b, lo, mid, b, mid+1, hi);
```

Problem.  Complex interactions with reverse copy.

Solution.  Go back to sentinels.

`Arrays.sort()`

## Sorting challenge 6

Problem.  Choose mergesort for Algs 4th edition.
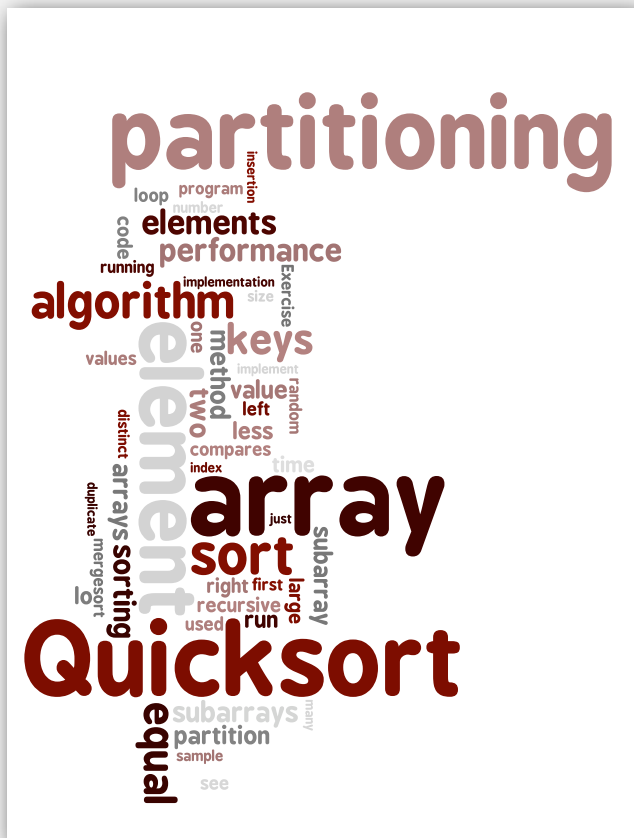
Recursive argument switch is out (recommended only for pros).

Q.  Why not use reverse array copy?

```
private static void merge(Comparable[] a, int lo, int mid, int hi)
{

   for (int i = lo; i <= mid; i++)
      aux[i] = a[i];


   for (int j = mid+1; j <= hi; j++)
      aux[j] = a[hi-j+mid+1];


   int i = lo, j = hi;
   for (int k = lo; k <= hi; k++)
      if (less(aux[j], aux[i])) a[k] = aux[j--];
      else                      a[k] = aux[i++];

}
```

# 2.3 Quicksort



‣ quicksort
‣ selection
‣ duplicate keys
‣ system sorts

## Two classic sorting algorithms

Critical components in the world's computational infrastructure.

- Full scientific understanding of their properties has enabled us to develop them into practical system sorts.
- Quicksort honored as one of top 10 algorithms of $20^{th}$ century in science and engineering.

### Mergesort.                    ⟵  last lecture

- Java sort for objects.
- Perl, Python stable sort.

### Quicksort.                    ⟵  this lecture

- Java sort for primitive types.
- C qsort, Unix, g++, Visual C++, Python.

# Quicksort

Basic plan.

- **Shuffle** the array.
- **Partition** so that, for some `j`
  - element `a[j]` is in place
  - no larger element to the left of `j`
  - no smaller element to the right of `j`
- **Sort** each piece recursively.

*Sir Charles Antony Richard Hoare*
*1980 Turing Award*

|  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **input** | Q | U | I | C | K | S | O | R | T | E | X | A | M | P | L | E |
| **shuffle** | K | R | A | T | E | L | E | P | U | I | M | Q | C | X | O | S |

*partitioning element*

| **partition** | E | C | A | I | E | K | L | P | U | T | M | Q | R | X | O | S |

*not greater*      *not less*

| **sort left** | A | C | E | E | I | K | L | P | U | T | M | Q | R | X | O | S |
| **sort right** | A | C | E | E | I | K | L | M | O | P | Q | R | S | T | U | X |
| **result** | A | C | E | E | I | K | L | M | O | P | Q | R | S | T | U | X |

**Quicksort overview**

4

# Quicksort partitioning

## Basic plan.

- Scan `i` from left for an item that belongs on the right.
- Scan `j` from right for item item that belongs on the left.
- Exchange `a[i]` and `a[j]`.
- Continue until pointers cross.

|  | i | j | v | a[i] |
|---|---|---|---|---|
|  |  |  | 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 |  |
| initial values | -1 | 15 | K R A T E L E P U I M Q C X O S |  |
| scan left, scan right | 1 | 12 | K R A T E L E P U I M Q C X O S |  |
| exchange | 1 | 12 | K C A T E L E P U I M Q R X O S |  |
| scan left, scan right | 3 | 9 | K C A T E L E P U I M Q R X O S |  |
| exchange | 3 | 9 | K C A I E L E P U T M Q R X O S |  |
| scan left, scan right | 5 | 6 | K C A I E L E P U T M Q R X O S |  |
| exchange | 5 | 6 | K C A I E E L P U T M Q R X O S |  |
| scan left, scan right | 6 | 5 | K C A I E E L P U T M Q R X O S |  |
| final exchange | 0 | 5 | E C A I E K L P U T M Q R X O S |  |
| result |  |  | E C A I E K L P U T M Q R X O S |  |

Partitioning trace (array contents before and after each exchange)

## Quicksort: Java code for partitioning

```java
private static int partition(Comparable[] a, int lo, int hi)
{
    int i = lo, j = hi+1;
    while (true)
    {
        while (less(a[++i], a[lo]))          find item on left to swap
            if (i == hi) break;

        while (less(a[lo], a[--j]))          find item on right to swap
            if (j == lo) break;

        if (i >= j) break;                   check if pointers cross
        exch(a, i, j);                       swap
    }

    exch(a, lo, j);                          swap with partitioning item
    return j;                   return index of item now known to be in place
}
```

| before | v | | | |
|--------|---|---|---|---|

↑ lo          ↑ hi

| during | v | ≤ v | | ≥ v |
|--------|---|-----|---|-----|

↑ i      ↑ j

| after | ≤ v | v | ≥ v |
|-------|-----|---|-----|

↑ lo     ↑ j        ↑ hi

6

# Quicksort: Java implementation

```java
public class Quick
{
   private static int partition(Comparable[] a, int lo, int hi)
   {  /* see previous slide */  }

   public static void sort(Comparable[] a)
   {
      StdRandom.shuffle(a);
      sort(a, 0, a.length - 1);
   }

   private static void sort(Comparable[] a, int lo, int hi)
   {
      if (hi <= lo) return;
      int j = partition(a, lo, hi);
      sort(a, lo, j-1);
      sort(a, j+1, hi);
   }
}
```

shuffle needed for performance guarantee

# Quicksort trace

| | lo | j | hi | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| initial values | | | | Q | U | I | C | K | S | O | R | T | E | X | A | M | P | L | E |
| random shuffle | | | | K | R | A | T | E | L | E | P | U | I | M | Q | C | X | O | S |
| | 0 | 5 | 15 | E | C | A | I | E | K | L | P | U | T | M | Q | R | X | O | S |
| | 0 | 3 | 4 | E | C | A | E | I | K | L | P | U | T | M | Q | R | X | O | S |
| | 0 | 2 | 2 | A | C | E | E | I | K | L | P | U | T | M | Q | R | X | O | S |
| | 0 | 0 | 1 | A | C | E | E | I | K | L | P | U | T | M | Q | R | X | O | S |
| | 1 | | 1 | A | C | E | E | I | K | L | P | U | T | M | Q | R | X | O | S |
| | 4 | | 4 | A | C | E | E | I | K | L | P | U | T | M | Q | R | X | O | S |
| | 6 | 6 | 15 | A | C | E | E | I | K | L | P | U | T | M | Q | R | X | O | S |
| | 7 | 9 | 15 | A | C | E | E | I | K | L | M | O | P | T | Q | R | X | U | S |
| | 7 | 7 | 8 | A | C | E | E | I | K | L | M | O | P | T | Q | R | X | U | S |
| | 8 | | 8 | A | C | E | E | I | K | L | M | O | P | T | Q | R | X | U | S |
| | 10 | 13 | 15 | A | C | E | E | I | K | L | M | O | P | S | Q | R | T | U | X |
| | 10 | 12 | 12 | A | C | E | E | I | K | L | M | O | P | R | Q | S | T | U | X |
| | 10 | 11 | 11 | A | C | E | E | I | K | L | M | O | P | Q | R | S | T | U | X |
| | 10 | | 10 | A | C | E | E | I | K | L | M | O | P | Q | R | S | T | U | X |
| | 14 | 14 | 15 | A | C | E | E | I | K | L | M | O | P | Q | R | S | T | U | X |
| | 15 | | 15 | A | C | E | E | I | K | L | M | O | P | Q | R | S | T | U | X |
| result | | | | A | C | E | E | I | K | L | M | O | P | Q | R | S | T | U | X |

*no partition for subarrays of size 1*

**Quicksort trace (array contents after each partition)**

# Quicksort animation

50 random elements

▲ algorithm position

in order

current subarray

not in order

## Quicksort: implementation details

**Partitioning in-place.** Using a spare array makes partitioning easier (and stable), but is not worth the cost.

**Terminating the loop.** Testing whether the pointers cross is a bit trickier than it might seem.

**Staying in bounds.** The `(j == lo)` test is redundant (why?), but the `(i == hi)` test is not.

**Preserving randomness.** Shuffling is needed for performance guarantee.

**Equal keys.** When duplicates are present, it is (counter-intuitively) best to stop on elements equal to the partitioning element.

## Quicksort: empirical analysis

Running time estimates:

- Home pc executes $10^8$ compares/second.

- Supercomputer executes $10^{12}$ compares/second.

| | insertion sort ($N^2$) | | | mergesort ($N \log N$) | | | quicksort ($N \log N$) | | |
|---|---|---|---|---|---|---|---|---|---|
| computer | thousand | million | billion | thousand | million | billion | thousand | million | billion |
| home | instant | 2.8 hours | 317 years | instant | 1 second | 18 min | instant | 0.3 sec | 6 min |
| super | instant | 1 second | 1 week | instant | instant | instant | instant | instant | instant |

Lesson 1. Good algorithms are better than supercomputers.

Lesson 2. Great algorithms are better than good ones.

# Quicksort: best case analysis

Best case.  Number of compares is ~ N lg N.

| lo | j | hi | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|----|---|----|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|
| | | | | | | | | | | a[ ] | | | | | | | |
| | | | H | A | C | B | F | E | G | D | L | I | K | J | N | M | O |
| | | | H | A | C | B | F | E | G | D | L | I | K | J | N | M | O |
| 0 | 7 | 14 | D | A | C | B | F | E | G | H | L | I | K | J | N | M | O |
| 0 | 3 | 6 | B | A | C | D | F | E | G | H | L | I | K | J | N | M | O |
| 0 | 1 | 2 | A | B | C | D | F | E | G | H | L | I | K | J | N | M | O |
| 0 | | 0 | A | B | C | D | F | E | G | H | L | I | K | J | N | M | O |
| 2 | | 2 | A | B | C | D | F | E | G | H | L | I | K | J | N | M | O |
| 4 | 5 | 6 | A | B | C | D | E | F | G | H | L | I | K | J | N | M | O |
| 4 | | 4 | A | B | C | D | E | F | G | H | L | I | K | J | N | M | O |
| 6 | | 6 | A | B | C | D | E | F | G | H | L | I | K | J | N | M | O |
| 8 | 11 | 14 | A | B | C | D | E | F | G | H | J | I | K | L | N | M | O |
| 8 | 9 | 10 | A | B | C | D | E | F | G | H | I | J | K | L | N | M | O |
| 8 | | 8 | A | B | C | D | E | F | G | H | I | J | K | L | N | M | O |
| 10 | | 10 | A | B | C | D | E | F | G | H | I | J | K | L | N | M | O |
| 12 | 13 | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 12 | | 12 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 14 | | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| | | | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |

# Quicksort: worst case analysis

Worst case. Number of compares is ~ $N^2 / 2$.

| lo | j | hi | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|----|----|----|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|
| | | | | | | | | | | a[ ] | | | | | | | |
| | | | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| | | | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 0 | 0 | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 1 | 1 | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 2 | 2 | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 3 | 3 | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 4 | 4 | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 5 | 5 | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 6 | 6 | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 7 | 7 | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 8 | 8 | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 9 | 9 | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 10 | 10 | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 11 | 11 | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 12 | 12 | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 13 | 13 | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 14 | | 14 | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| | | | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |

## Quicksort: average-case analysis

Proposition I. The average number of compares $C_N$ to quicksort an array of N elements is ~ 2N ln N (and the number of exchanges is ~ ⅓ N ln N).

Pf. $C_N$ satisfies the recurrence $C_0 = C_1 = 0$ and for N ≥ 2:

$$C_N = (N+1) + \frac{C_0 + C_1 + \ldots + C_{N-1}}{N} + \frac{C_{N-1} + C_{N-2} + \ldots + C_0}{N}$$

↑ partitioning     ↑ left     ↑ right     ↖ partitioning probability

- Multiply both sides by N and collect terms:

$$NC_N = N(N+1) + 2(C_0 + C_1 + \ldots + C_{N-1})$$

- Subtract this from the same equation for N-1:

$$NC_N - (N-1)C_{N-1} = 2N + 2C_{N-1}$$

- Rearrange terms and divide by N(N+1):

$$\frac{C_N}{N+1} = \frac{C_{N-1}}{N} + \frac{2}{N+1}$$

## Quicksort: average-case analysis

- Repeatedly apply above equation:

$$\frac{C_N}{N+1} = \frac{C_{N-1}}{N} + \frac{2}{N+1}$$

$$= \frac{C_{N-2}}{N-1} + \frac{2}{N} + \frac{2}{N+1}$$

$$= \frac{C_{N-3}}{N-2} + \frac{2}{N-1} + \frac{2}{N} + \frac{2}{N+1}$$

$$= \frac{2}{1} + \frac{2}{2} + \frac{2}{3} + \ldots + \frac{2}{N+1}$$

*previous equation*

- Approximate sum by an integral:

$$C_N \sim 2(N+1)\left(1 + \frac{1}{2} + \frac{1}{3} + \ldots \frac{1}{N}\right)$$

$$\sim 2(N+1)\int_1^N \frac{1}{x}dx$$



- Finally, the desired result:

$$C_N \sim 2(N+1)\ln N \approx 1.39 N \lg N$$

## Quicksort: summary of performance characteristics

**Worst case.** Number of compares is quadratic.

- N + (N-1) + (N-2) + … + 1  ~ $N^2$ / 2.
- More likely that your computer is struck by lightning.

**Average case.** Number of compares is ~ 1.39 N lg N.

- 39% more compares than mergesort.
- But faster than mergesort in practice because of less data movement.

**Random shuffle.**

- Probabilistic guarantee against worst case.
- Basis for math model that can be validated with experiments.

**Caveat emptor.** Many textbook implementations go quadratic if input:

- Is sorted or reverse sorted.
- Has many duplicates (even if randomized!)  [stay tuned]

# Quicksort: practical improvements

## Median of sample.

- Best choice of pivot element = median.
- Estimate true median by taking median of sample.

## Insertion sort small subarrays.

- Even quicksort has too much overhead for tiny subarrays.
- Can delay insertion sort until end.

## Optimize parameters.

~ 12/7  N ln N compares
~ 12/35 N ln N exchanges

- Median-of-3 random elements.
- Cutoff to insertion sort for ≈ 10 elements.

## Non-recursive version.

guarantees O(log N) stack size

- Use explicit stack.
- Always sort smaller half first.

# Quicksort with cutoff to insertion sort:  visualization



Quicksort with median-of-3 partitioning and cutoff for small subarrays

‣ quicksort
‣ **selection**
‣ duplicate keys
‣ system sorts

## Selection

Goal. Find the $k^{th}$ largest element.

Ex. Min (k = 0), max (k = N-1), median (k = N/2).

Applications.
- Order statistics.
- Find the "top k."

Use theory as a guide.
- Easy $O(N \log N)$ upper bound.
- Easy $O(N)$ upper bound for k = 1, 2, 3.
- Easy $\Omega(N)$ lower bound.

Which is true?
- $\Omega(N \log N)$ lower bound? ⟵ is selection as hard as sorting?
- $O(N)$ upper bound? ⟵ is there a linear-time algorithm for all k?

## Quick-select

Partition array so that:

- Element `a[j]` is in place.
- No larger element to the left of `j`.
- No smaller element to the right of `j`.

Repeat in one subarray, depending on `j`; finished when `j` equals `k`.

```java
public static Comparable select(Comparable[] a, int k)
{
    StdRandom.shuffle(a);
    int lo = 0, hi = a.length - 1;
    while (hi > lo)
    {
        int j = partition(a, lo, hi);
        if        (j < k) lo = j + 1;
        else if (j > k) hi = j - 1;
        else              return a[k];
    }
    return a[k];
}
```

if a[k] is here
set hi to j−1

if a[k] is here
set lo to j+1

| ≤ v | v | ≥ v |

lo          j          hi

## Quick-select:  mathematical analysis

Proposition.  Quick-select takes linear time on average.

Pf sketch.
- Intuitively, each partitioning step roughly splits array in half:

  N + N/2 + N/4 + … + 1  ~ 2N compares.
- Formal analysis similar to quicksort analysis yields:

$$C_N = 2N + k \ln(N/k) + (N-k) \ln(N/(N-k))$$

Ex.  (2 + 2 ln 2) N compares to find the median.

Remark.  Quick-select uses $\sim N^2/2$ compares in worst case,
but as with quicksort, the random shuffle provides a probabilistic guarantee.

## Theoretical context for selection

Challenge.  Design algorithm whose worst-case running time is linear.

Proposition.  [Blum, Floyd, Pratt, Rivest, Tarjan, 1973] There exists a compare-based selection algorithm whose worst-case running time is linear.

Remark.  But, algorithm is too complicated to be useful in practice.

Use theory as a guide.
- Still worthwhile to seek practical linear-time (worst-case) algorithm.
- Until one is discovered, use quick-select if you don't need a full sort.

## Generic methods

In our `select()` implementation, client needs a cast.

```
Double[] a = new Double[N];
for (int i = 0; i < N; i++)
    a[i] = StdRandom.uniform();
Double median = (Double) Quick.select(a, N/2);
```

unsafe cast
required

The compiler also complains.

```
% javac Quick.java
Note: Quick.java uses unchecked or unsafe operations.
Note: Recompile with -Xlint:unchecked for details.
```

Q. How to fix?

# Generic methods

Pedantic (safe) version.  Compiles cleanly, no cast needed in client.

```
public class QuickPedantic                    generic type variable
{                                             (value inferred from argument a[])
    public  static <Key extends Comparable<Key>> Key select(Key[] a, int k)
    {  /* as before */  }
                                                              return type matches array type

    public  static <Key extends Comparable<Key>> void sort(Key[] a)
    {  /* as before */  }


    private static <Key extends Comparable<Key>> int partition(Key[] a, int lo, int hi)
    {  /* as before */  }


    private static <Key extends Comparable<Key>> boolean less(Key v, Key w)
    {  /* as before */  }


    private static <Key extends Comparable<Key>> void exch(Key[] a, int i, int j)
    {  Key swap = a[i]; a[i] = a[j]; a[j] = swap;  }

}                 can declare variables of generic type
```

http://www.cs.princeton.edu/algs4/35applications/QuickPedantic.java.html

Remark.  Obnoxious code needed in system sort; not in this course (for brevity).

- quicksort
- selection
- **duplicate keys**
- system sorts

## Duplicate keys

Often, purpose of sort is to bring records with duplicate keys together.

- Sort population by age.
- Find collinear points.  ← see Assignment 3
- Remove duplicates from mailing list.
- Sort job applicants by college attended.

Typical characteristics of such applications.

- Huge array.
- Small number of key values.

```
Chicago 09:25:52
Chicago 09:03:13
Chicago 09:21:05
Chicago 09:19:46
Chicago 09:19:32
Chicago 09:00:00
Chicago 09:35:21
Chicago 09:00:59
Houston 09:01:10
Houston 09:00:13
Phoenix 09:37:44
Phoenix 09:00:03
Phoenix 09:14:25
Seattle 09:10:25
Seattle 09:36:14
Seattle 09:22:43
Seattle 09:10:11
Seattle 09:22:54
```

key

## Duplicate keys

Mergesort with duplicate keys.  Always ~ N lg N compares.

Quicksort with duplicate keys.

- Algorithm goes quadratic unless partitioning stops on equal keys!
- 1990s C user found this defect in `qsort()`.

several textbook and system implementations
also have this defect

**S T O P O N E Q U A L K E Y S**

swap

if we don't stop
on equal keys

if we stop on
equal keys

# Duplicate keys:  the problem

**Mistake.** Put all keys equal to the partitioning element on one side.

**Consequence.**   ~ $N^2 / 2$ compares when all keys equal.

        B  A  A  B  A  B  B  **B**  C  C  C        A  A  A  A  A  A  A  A  A  A  **A**

**Recommended.**  Stop scans on keys equal to the partitioning element.

**Consequence.**  ~ $N \lg N$ compares when all keys equal.

        B  A  A  B  A  **B**  C  C  B  C  B        A  A  A  A  A  **A**  A  A  A  A

**Desirable.**  Put all keys equal to the partitioning element in place.

        A  A  A  **B  B  B  B  B**  C  C  C        **A  A  A  A  A  A  A  A  A  A**

## 3-way partitioning

Goal. Partition array into 3 parts so that:

- Elements between `lt` and `gt` equal to partition element `v`.
- No larger elements to left of `lt`.
- No smaller elements to right of `gt`.



Dutch national flag problem. [Edsger Dijkstra]

- Conventional wisdom until mid 1990s: not worth doing.
- New approach discovered when fixing mistake in C library `qsort()`.
- Now incorporated into `qsort()` and Java system sort.

## 3-way partitioning:  Dijkstra's solution

3-way partitioning.

- Let `v` be partitioning element `a[lo]`.
- Scan `i` from left to right.
  - `a[i]` less than `v`: exchange `a[lt]` with `a[i]` and increment both `lt` and `i`
  - `a[i]` greater than `v`: exchange `a[gt]` with `a[i]` and decrement `gt`
  - `a[i]` equal to `v`: increment `i`

All the right properties.

- In-place.
- Not much code.
- Small overhead if no equal keys.

**3-way partitioning**

# 3-way partitioning:  trace

| lt | i | gt | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|----|---|----|---|---|---|---|---|---|---|---|---|---|---|----|----|
| | | | a[] | | | | | | | | | | | | |
| 0 | 0 | 11 | | R | B | W | W | R | W | B | R | R | W | B | R |
| 0 | 1 | 11 | | R | B | W | W | R | W | B | R | R | W | B | R |
| 1 | 2 | 11 | | B | R | W | W | R | W | B | R | R | W | B | R |
| 1 | 2 | 10 | | B | R | R | W | R | W | B | R | R | W | B | W |
| 1 | 3 | 10 | | B | R | R | W | R | W | B | R | R | W | B | W |
| 1 | 3 | 9 | | B | R | R | B | R | W | B | R | R | W | W | W |
| 2 | 4 | 9 | | B | B | R | R | R | W | B | R | R | W | W | W |
| 2 | 5 | 9 | | B | B | R | R | R | W | B | R | R | W | W | W |
| 2 | 5 | 8 | | B | B | R | R | R | W | B | R | R | W | W | W |
| 2 | 5 | 7 | | B | B | R | R | R | R | B | R | W | W | W | W |
| 2 | 6 | 7 | | B | B | R | R | R | R | B | R | W | W | W | W |
| 3 | 7 | 7 | | B | B | B | R | R | R | R | R | W | W | W | W |
| 3 | 8 | 7 | | B | B | B | R | R | R | R | R | W | W | W | W |

**3-way partitioning trace (array contents after each loop iteration)**

## 3-way quicksort: Java implementation

```java
private static void sort(Comparable[] a, int lo, int hi)
{
   if (hi <= lo) return;
   int lt = lo, gt = hi;
   Comparable v = a[lo];
   int i = lo;
   while (i <= gt)
   {
      int cmp = a[i].compareTo(v);
      if      (cmp < 0) exch(a, lt++, i++);
      else if (cmp > 0) exch(a, i, gt--);
      else              i++;
   }

   sort(a, lo, lt - 1);
   sort(a, gt + 1, hi);
}
```



before  v

during  <v   =v          >v

after   <v       =v        >v

# 3-way quicksort: visual trace



*equal to partitioning element*

**Visual trace of quicksort with 3-way partitioning**

# Duplicate keys:  lower bound

Sorting lower bound.  If there are n distinct keys and the $i^{th}$ one occurs
$x_i$ times, any compare-based sorting algorithm must use at least

$$\lg \left( \frac{N!}{x_1! \; x_2! \; \cdots \; x_n!} \right) \;\; \sim \;\; - \sum_{i=1}^{n} x_i \lg \frac{x_i}{N}$$

<span style="color:red">N lg N when all distinct;<br>linear when only a constant number of distinct keys</span>

compares in the worst case.

<span style="color:red">proportional to lower bound</span>

Proposition.  [Sedgewick-Bentley, 1997]
Quicksort with 3-way partitioning is entropy-optimal.
Pf.  [beyond scope of course]

Bottom line.  Randomized quicksort with 3-way partitioning reduces running
time from linearithmic to linear in broad class of applications.

- selection
- duplicate keys
- comparators
- **system sorts**

## Sorting applications

### Sorting algorithms are essential in a broad variety of applications:

- Sort a list of names.
- Organize an MP3 library.
- Display Google PageRank results.
- List RSS news items in reverse chronological order.

obvious applications

- Find the median.
- Find the closest pair.
- Binary search in a database.
- Identify statistical outliers.
- Find duplicates in a mailing list.

problems become easy once items
are in sorted order

- Data compression.
- Computer graphics.
- Computational biology.
- Supply chain management.
- Load balancing on a parallel computer.

. . .

non-obvious applications

### Every system needs (and has) a system sort!

## Java system sorts

Java uses both mergesort and quicksort.

- **Arrays.sort()** sorts array of **Comparable** or any primitive type.
- Uses quicksort for primitive types; mergesort for objects.

```java
import java.util.Arrays;

public class StringSort
{
    public static void main(String[] args)
    {
        String[] a = StdIn.readAll().split("\\s+");
        Arrays.sort(a);
        for (int i = 0; i < N; i++)
            StdOut.println(a[i]);
    }
}
```

Q. Why use different algorithms, depending on type?

## Java system sort for primitive types

**Engineering a sort function.** [Bentley-McIlroy, 1993]

- Original motivation: improve `qsort()`.
- Basic algorithm = 3-way quicksort with cutoff to insertion sort.
- Partition on Tukey's ninther: median of the medians of 3 samples, each of 3 elements.

approximate median-of-9

| nine evenly spaced elements | R | L | A | P | M | C | G | A | X | Z | K | R | B | R | J | J | E |

| groups of 3 | R | A | M | | G | X | K | | B | J | E |

| medians | M | K | E |

| ninther | K |

## Why use Tukey's ninther?

- Better partitioning than random shuffle.
- Less costly than random shuffle.

## Achilles heel in Bentley-McIlroy implementation (Java system sort)

Based on all this research, Java's system sort is solid, right?

A killer input.

more disastrous consequences in C

- Blows function call stack in Java and crashes program.
- Would take quadratic time if it didn't crash first.

```
% more 250000.txt
0
218750
222662
11
166672
247070
83339
...
```

```
% java IntegerSort < 250000.txt
Exception in thread "main"
java.lang.StackOverflowError
    at java.util.Arrays.sort1(Arrays.java:562)
    at java.util.Arrays.sort1(Arrays.java:606)
    at java.util.Arrays.sort1(Arrays.java:608)
    at java.util.Arrays.sort1(Arrays.java:608)
    at java.util.Arrays.sort1(Arrays.java:608)
    ...
```

250,000 integers
between 0 and 250,000

Java's sorting library crashes, even if
you give it as much stack space as Windows allows

McIlroy's devious idea. [A Killer Adversary for Quicksort]

- Construct malicious input while running system quicksort,
  in response to elements compared.
- If `v` is partitioning element, commit to `(v < a[i])` and `(v < a[j])`, but don't
  commit to `(a[i] < a[j])` or `(a[j] > a[i])` until `a[i]` and `a[j]` are compared.

Consequences.

- Confirms theoretical possibility.
- Algorithmic complexity attack:  you enter linear amount of data;
  server performs quadratic amount of work.

Remark.  Attack is not effective if array is shuffled before sort.

Q.  Why do you think system sort is deterministic?

## System sort: Which algorithm to use?

Many sorting algorithms to choose from:

### Internal sorts.
- Insertion sort, selection sort, bubblesort, shaker sort.
- Quicksort, mergesort, heapsort, samplesort, shellsort.
- Solitaire sort, red-black sort, splaysort, Dobosiewicz sort, psort, ...

### External sorts.  Poly-phase mergesort, cascade-merge, oscillating sort.

### Radix sorts.  Distribution, MSD, LSD, 3-way radix quicksort.

### Parallel sorts.
- Bitonic sort, Batcher even-odd sort.
- Smooth sort, cube sort, column sort.
- GPUsort.

## System sort: Which algorithm to use?

Applications have diverse attributes.

- Stable?
- Parallel?
- Deterministic?
- Keys all distinct?
- Multiple key types?
- Linked list or arrays?
- Large or small records?
- Is your array randomly ordered?
- Need guaranteed performance?



many more combinations of attributes than algorithms

Elementary sort may be method of choice for some combination.

Cannot cover all combinations of attributes.

Q.  Is the system sort good enough?

A.  Usually.

## Sorting summary

| | inplace? | stable? | worst | average | best | remarks |
|---|:---:|:---:|:---:|:---:|:---:|:---:|
| selection | × | | $N^2/2$ | $N^2/2$ | $N^2/2$ | $N$ exchanges |
| insertion | × | × | $N^2/2$ | $N^2/4$ | $N$ | use for small $N$ or partially ordered |
| shell | × | | ? | ? | $N$ | tight code, subquadratic |
| quick | × | | $N^2/2$ | $2\,N \ln N$ | $N \lg N$ | $N \log N$ probabilistic guarantee fastest in practice |
| 3-way quick | × | | $N^2/2$ | $2\,N \ln N$ | $N$ | improves quicksort in presence of duplicate keys |
| merge | | × | $N \lg N$ | $N \lg N$ | $N \lg N$ | $N \log N$ guarantee, stable |
| ??? | × | × | $N \lg N$ | $N \lg N$ | $N \lg N$ | holy sorting grail |

# Which sorting algorithm?

| lifo | find | data | data | data | data | hash | data |
|------|------|------|------|------|------|------|------|
| fifo | fifo | fifo | fifo | exch | fifo | fifo | exch |
| data | data | find | find | fifo | lifo | data | fifo |
| type | exch | hash | hash | find | type | link | find |
| hash | hash | heap | heap | hash | hash | leaf | hash |
| heap | heap | lifo | lifo | heap | heap | heap | heap |
| sort | less | link | link | leaf | link | exch | leaf |
| link | left | list | list | left | sort | node | left |
| list | leaf | push | push | less | find | lifo | less |
| push | lifo | root | root | lifo | list | left | lifo |
| find | push | sort | sort | link | push | find | link |
| root | root | type | type | list | root | path | list |
| leaf | list | leaf | leaf | sort | leaf | list | next |
| tree | tree | left | tree | tree | null | next | node |
| null | null | node | null | null | path | less | null |
| path | path | null | path | path | tree | root | path |
| node | node | path | node | node | exch | sink | push |
| left | link | tree | left | type | left | swim | root |
| less | sort | exch | less | root | less | null | sink |
| exch | type | less | exch | push | node | sort | sort |
| sink | sink | next | sink | sink | next | type | swap |
| swim | swim | sink | swim | swim | sink | tree | swim |
| next | next | swap | next | next | swap | push | tree |
| swap | swap | swim | swap | swap | swim | swap | type |
| original | ? | ? | ? | ? | ? | ? | sorted |

# 2.4 Priority Queues

▸ API
▸ elementary implementations
▸ binary heaps
▸ heapsort
▸ event-based simulation

---

## Priority queue API

| data type | delete |
|---|---|
| stack | last in, first out |
| queue | first in, first out |
| priority queue | largest value out |

```
public class MaxPQ<Key extends Comparable<Key>>

        MaxPQ()              create a priority queue

        MaxPQ(maxN)          create a priority queue of initial capacity maxN

  void  insert(Key v)        insert a key into the priority queue

   Key  max()                return the largest key

   Key  delMax()             return and remove the largest key

boolean isEmpty()            is the priority queue empty?

   int  size()               number of entries in the priority queue
```

**API for a generic priority queue**

| operation | argument | return value |
|---|---|---|
| insert | P | |
| insert | Q | |
| insert | E | |
| remove max | | Q |
| insert | X | |
| insert | A | |
| insert | M | |
| remove max | | X |
| insert | P | |
| insert | L | |
| insert | E | |
| remove max | | P |

2

---

## Priority queue applications

- Event-driven simulation.          [customers in a line, colliding particles]
- Numerical computation.            [reducing roundoff error]
- Data compression.                 [Huffman codes]
- Graph searching.                  [Dijkstra's algorithm, Prim's algorithm]
- Computational number theory.      [sum of powers]
- Artificial intelligence.          [A* search]
- Statistics.                       [maintain largest M values in a sequence]
- Operating systems.                [load balancing, interrupt handling]
- Discrete optimization.            [bin packing, scheduling]
- Spam filtering.                   [Bayesian spam filter]

**Generalizes:** stack, queue, randomized queue.

3

---

## Priority queue client example

**Problem.** Find the largest M in a stream of N elements.
- Fraud detection: isolate $$ transactions.
- File maintenance: find biggest files or directories.

**Constraint.** Not enough memory to store N elements.

**Solution.** Use a min-oriented priority queue.

```
MinPQ<String> pq = new MinPQ<String>();

while(!StdIn.isEmpty())
{
   String s = StdIn.readString();
   pq.insert(s);
   if (pq.size() > M)
      pq.delMin();
}

while (!pq.isEmpty())
   System.out.println(pq.delMin());
```

| implementation | time | space |
|---|---|---|
| sort | N log N | N |
| elementary PQ | M N | M |
| binary heap | N log M | M |
| best in theory | N | M |

cost of finding the largest M
in a stream of N items

4

## Slide 5

5

## Slide 6

Priority queue:  unordered and ordered array implementation

| operation | argument | return value | size | contents (unordered) | | | | | | contents (ordered) | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| insert | P | | 1 | P | | | | | | P | | | | | |
| insert | Q | | 2 | P | Q | | | | | P | Q | | | | |
| insert | E | | 3 | P | Q | E | | | | E | P | Q | | | |
| remove max | | Q | 2 | P | E | | | | | E | P | | | | |
| insert | X | | 3 | P | E | X | | | | E | P | X | | | |
| insert | A | | 4 | P | E | X | A | | | A | E | P | X | | |
| insert | M | | 5 | P | E | X | A | M | | A | E | M | P | X | |
| remove max | | X | 4 | P | E | M | A | | | A | E | M | P | | |
| insert | P | | 5 | P | E | M | A | P | | A | E | M | P | P | |
| insert | L | | 6 | P | E | M | A | P | L | A | E | L | M | P | P |
| insert | E | | 7 | P | E | M | A | P | L | E | A | E | E | L | M | P | P |
| remove max | | P | 6 | E | M | A | P | L | E | A | E | E | L | M | P |

**A sequence of operations on a priority queue**

6

## Slide 7

Priority queue:  unordered array implementation

```
public class UnorderedMaxPQ<Key extends Comparable<Key>>
{
    private Key[] pq;    // pq[i] = ith element on pq
    private int N;       // number of elements on pq

    public UnorderedMaxPQ(int capacity)
    {  pq = (Key[]) new Comparable[capacity];  }

    public boolean isEmpty()
    {  return N == 0;  }

    public void insert(Key x)
    {  pq[N++] = x;  }

    public Key delMax()
    {
        int max = 0;
        for (int i = 1; i < N; i++)
            if (less(max, i)) max = i;
        exch(max, N-1);
        return pq[--N];
    }
}
```

no generic array creation

less() and exch() as for sorting

7

## Slide 8

Priority queue elementary implementations

Challenge.  Implement all operations efficiently.

| implementation | insert | del max | max |
|---|---|---|---|
| unordered array | 1 | N | N |
| ordered array | N | 1 | 1 |
| goal | log N | log N | log N |

order-of-growth running time for PQ with N items

8

## Slide 9

## Slide 10

Binary tree

Binary tree. Empty or node with links to left and right binary trees.

Complete tree. Perfectly balanced, except for bottom level.



complete tree of height 5

N = 16
⌊lg N⌋ = 4
height = 5

Property. Height of complete tree with N nodes is 1 + ⌊lg N⌋.
Pf. Height only increases when N is exactly a power of 2.

## Slide 11

A complete binary tree in nature



Hyphaene Compressa - Doum Palm        © Shlomit Pinter

## Slide 12

Binary heap

Binary heap. Array representation of a heap-ordered complete binary tree.

Heap-ordered binary tree.
• Keys in nodes.
• No smaller than children's keys.

Array representation.
• Take nodes in level order.
• No explicit links needed!

| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|---|----|----|
| a[i] | – | T | S | R | P | N | O | A | E | I | H | G |

## Binary heap properties

**Property A.** Largest key is `a[1]`, which is root of binary tree.

indices start at 1

**Property B.** Can use array indices to move through tree.
- Parent of node at `k` is at `k/2`.
- Children of node at `k` are at `2k` and `2k+1`.

---

## Promotion in a heap

**Scenario.** Node's key becomes larger key than its parent's key.

**To eliminate the violation:**
- Exchange key in node with key in parent.
- Repeat until heap order restored.

```
private void swim(int k)
{
    while (k > 1 && less(k/2, k))
    {
        exch(k, k/2);
        k = k/2;
    }
}
```

parent of node at k is at k/2



violates heap order
(larger key than parent)

**Peter principle.** Node promoted to level of incompetence.

---

## Insertion in a heap

**Insert.** Add node at end, then swim it up.
**Running time.** At most ~ lg N compares.

```
public void insert(Key x)
{
    pq[++N] = x;
    swim(N);
}
```

insert

S ← key to insert

add key to heap
violates heap order

swim up

---

## Demotion in a heap

**Scenario.** Node's key becomes smaller than one (or both) of its children's keys.

**To eliminate the violation:**
- Exchange key in node with key in larger child.
- Repeat until heap order restored.

```
private void sink(int k)
{
    while (2*k <= N)
    {
        int j = 2*k;
        if (j < N && less(j, j+1)) j++;
        if (!less(k, j)) break;
        exch(k, j);
        k = j;
    }
}
```

children of node
at k are 2k and 2k+1

violates heap order
(smaller than a child)



**Power struggle.** Better subordinate promoted.

## Delete the maximum in a heap

**Delete max.** Exchange root with node at end, then sink it down.

**Running time.** At most ~ 2 lg N compares.

```
public Key delMax()
{
   Key max = pq[1];
   exch(1, N--);
   sink(1);
   pq[N+1] = null;      ← prevent loitering
   return max;
}
```



remove the maximum

T ← key to remove

exchange keys with root

H ← violates heap order

T ← remove node from heap

sink down

---

## Heap operations



insert P

insert Q

insert E

remove max (Q)

insert X

insert A

insert M

remove max (X)

insert P

insert L

insert E

remove max (P)

---

## Binary heap:  Java implementation

```
public class MaxPQ<Key extends Comparable<Key>>
{
   private Key[] pq;
   private int N;

   public MaxPQ(int capacity)
   {  pq = (Key[]) new Comparable[capacity+1];  }

   public boolean isEmpty()
   {    return N == 0;   }
   public void insert(Key key)
   {   /* see previous code */  }       ← PQ ops
   public Key delMax()
   {   /* see previous code */  }

   private void swim(int k)
   {   /* see previous code */  }       ← heap helper functions
   private void sink(int k)
   {   /* see previous code */  }

   private boolean less(int i, int j)
   {   return pq[i].compareTo(pq[j]) < 0;  }    ← array helper functions
   private void exch(int i, int j)
   {   Key t = pq[i]; pq[i] = pq[j]; pq[j] = t;  }
}
```

---

## Priority queues implementation cost summary

| implementation | insert | del max | max |
|---|---|---|---|
| unordered array | 1 | N | N |
| ordered array | N | 1 | 1 |
| binary heap | log N | log N | 1 |

order-of-growth running time for PQ with N items

**Hopeless challenge.**  Make all operations constant time.

**Q.**  Why hopeless?

## Binary heap considerations

Minimum-oriented priority queue.
- Replace `less()` with `greater()`.
- Implement `greater()`.

Dynamic array resizing.
- Add no-arg constructor.
- Apply repeated doubling and shrinking.  ⟵ leads to O(log N) amortized time per op

Immutability of keys.
- Assumption: client does not change keys while they're on the PQ.
- Best practice: use immutable keys.

Other operations.
- Remove an arbitrary item.  ⟵
- Change the priority of an item.  ⟵ easy to implement with `sink()` and `swim()` [stay tuned]

‣ API
‣ elementary implementations
‣ binary heaps
‣ **heapsort**
‣ event-based simulation

## Heapsort

Basic plan for in-place sort.
- Create max-heap with all N keys.
- Repeatedly remove the maximum key.



start with array of keys in arbitrary order

build a max-heap (in place)

sorted result (in place)

## Heapsort: heap construction

First pass. Build heap using bottom-up method.

```
for (int k = N/2; k >= 1; k--)
    sink(a, k, N);
```

## Heapsort: sortdown

Second pass.

- Remove the maximum, one at a time.
- Leave in array, instead of nulling out.

```
while (N > 1)
{
    exch(a, 1, N--);
    sink(a, 1, N);
}
```

---

## Heapsort: Java implementation

```
public class Heap
{
   public static void sort(Comparable[] pq)
   {
      int N = pq.length;
      for (int k = N/2; k >= 1; k--)
         sink(pq, k, N);
      while (N > 1)
      {
         exch(pq, 1, N);
         sink(pq, 1, --N);
      }
   }

   private static void sink(Comparable[] pq, int k, int N)
   {  /* as before */  }

   private static boolean less(Comparable[] pq, int i, int j)
   {  /* as before */  }

   private static void exch(Comparable[] pq, int i, int j)
   {  /* as before */  }

}
```

but use 1-based indexing

---

## Heapsort: trace



**Heapsort trace (array contents just after each sink)**

---

## Heapsort: mathematical analysis

Proposition Q. At most 2 N lg N compares and exchanges.

Significance. Sort in N log N worst-case without using extra memory.

- Mergesort: no, linear extra space. ⟵ in-place merge possible, not practical
- Quicksort: no, quadratic time in worst case. ⟵ N log N worst-case quicksort possible, not practical
- Heapsort: yes!

Bottom line. Heapsort is optimal for both time and space, but:

- Inner loop longer than quicksort's.
- Makes poor use of cache memory.
- Not stable.

50 random elements



▲ algorithm position

in order

not in order

http://www.sorting-algorithms.com/heap-sort

|  | inplace? | stable? | worst | average | best | remarks |
|---|---|---|---|---|---|---|
| selection | × |  | $N^2/2$ | $N^2/2$ | $N^2/2$ | $N$ exchanges |
| insertion | × | × | $N^2/2$ | $N^2/4$ | $N$ | use for small $N$ or partially ordered |
| shell | × |  | ? | ? | $N$ | tight code, subquadratic |
| quick | × |  | $N^2/2$ | $2N \ln N$ | $N \lg N$ | $N \log N$ probabilistic guarantee fastest in practice |
| 3-way quick | × |  | $N^2/2$ | $2N \ln N$ | $N$ | improves quicksort in presence of duplicate keys |
| merge |  | × | $N \lg N$ | $N \lg N$ | $N \lg N$ | $N \log N$ guarantee, stable |
| heap | × |  | $2N \lg N$ | $2N \lg N$ | $N \lg N$ | $N \log N$ guarantee, in-place |
| ??? | × | × | $N \lg N$ | $N \lg N$ | $N \lg N$ | holy sorting grail |

‣ API
‣ elementary implementations
‣ binary heaps
‣ heapsort
‣ **event-based simulation**

Goal. Simulate the motion of N moving particles that behave according to the laws of elastic collision.

## Molecular dynamics simulation of hard discs

**Goal.** Simulate the motion of N moving particles that behave according to the laws of elastic collision.

**Hard disc model.**
- Moving particles interact via elastic collisions with each other and walls.
- Each particle is a disc with known position, velocity, mass, and radius.
- No other forces.

*temperature, pressure, diffusion constant*   *motion of individual atoms and molecules*

**Significance.** Relates macroscopic observables to microscopic dynamics.
- Maxwell-Boltzmann: distribution of speeds as a function of temperature.
- Einstein: explain Brownian motion of pollen grains.

## Warmup: bouncing balls

**Time-driven simulation.** N bouncing balls in the unit square.

```
public class BouncingBalls
{
   public static void main(String[] args)
   {
      int N = Integer.parseInt(args[0]);
      Ball balls[] = new Ball[N];
      for (int i = 0; i < N; i++)
         balls[i] = new Ball();
      while(true)
      {
         StdDraw.clear();
         for (int i = 0; i < N; i++)
         {
            balls[i].move(0.5);
            balls[i].draw();
         }
         StdDraw.show(50);
      }
   }
}
```

`% java BouncingBalls 100`

*main simulation loop*

## Warmup: bouncing balls

```
public class Ball
{
    private double rx, ry;        // position
    private double vx, vy;        // velocity
    private final double radius;  // radius
    public Ball()
    { /* initialize position and velocity */  }

    public void move(double dt)
    {
        if ((rx + vx*dt < radius) || (rx + vx*dt > 1.0 - radius)) { vx = -vx; }
        if ((ry + vy*dt < radius) || (ry + vy*dt > 1.0 - radius)) { vy = -vy; }
        rx = rx + vx*dt;
        ry = ry + vy*dt;
    }
    public void draw()
    { StdDraw.filledCircle(rx, ry, radius);  }
}
```

*check for collision with walls*

**Missing.** Check for balls colliding with each other.
- Physics problems: when? what effect?
- CS problems: which object does the check? too many checks?

## Time-driven simulation

- Discretize time in quanta of size dt.
- Update the position of each particle after every dt units of time, and check for overlaps.
- If overlap, roll back the clock to the time of the collision, update the velocities of the colliding particles, and continue the simulation.

t

t + dt

t + 2 dt
(collision detected)

t + Δt
(roll back clock)

## Time-driven simulation

Main drawbacks.
- ~ $N^2/2$ overlap checks per time quantum.
- Simulation is too slow if dt is very small.
- May miss collisions if dt is too large.
  (if colliding particles fail to overlap when we are looking)



**dt too small: excessive computation**

**dt too large: may miss collisions**

## Event-driven simulation

Change state only when something happens.
- Between collisions, particles move in straight-line trajectories.
- Focus only on times when collisions occur.
- Maintain PQ of collision events, prioritized by time.
- Remove the min = get next collision.

Collision prediction. Given position, velocity, and radius of a particle, when will it collide next with a wall or another particle?

Collision resolution. If collision occurs, update colliding particle(s) according to laws of elastic collisions.



**prediction (at time $t$)**
*particles hit unless one passes intersection point before the other arrives (see Exercise 3.6.X)*

**resolution (at time $t + dt$)**
*velocities of both particles change after collision*

## Particle-wall collision

Collision prediction and resolution.
- Particle of radius $s$ at position $(rx, ry)$.
- Particle moving in unit box with velocity $(vx, vy)$.
- Will it collide with a vertical wall? If so, when?



**resolution (at time $t + dt$)**
*velocity after collision  $= ( - v_x, v_y)$*
*position after collision  $= ( 1 - s, r_y + v_y dt)$*

**prediction (at time $t$)**
*$dt \equiv$ time to hit wall*
*$= distance/velocity$*
*$= (1 - s - r_x)/v_x$*

$(r_x, r_y)$

$v_x$   $v_y$

$1 - s - r_x$

*wall at*
*$x = 1$*

$s$

**Predicting and resolving a particle-wall collision**

## Particle-particle collision prediction

Collision prediction.
- Particle $i$:  radius $s_i$, position $(rx_i, ry_i)$, velocity $(vx_i, vy_i)$.
- Particle $j$:  radius $s_j$, position $(rx_j, ry_j)$, velocity $(vx_j, vy_j)$.
- Will particles $i$ and $j$ collide? If so, when?



$(vx_i', vy_i')$

$(vx_j', vy_j')$

$m_i$

$s_i$

$(vx_i, vy_i)$

$(rx_i, ry_i)$

$(rx_i', ry_i')$

$i$

*time = $t$*

*time = $t + \Delta t$*

$(vx_j, vy_j)$

$s_j$

$j$

## Particle-particle collision prediction

Collision prediction.

- Particle $i$: radius $s_i$, position $(rx_i, ry_i)$, velocity $(vx_i, vy_i)$.
- Particle $j$: radius $s_j$, position $(rx_j, ry_j)$, velocity $(vx_j, vy_j)$.
- Will particles $i$ and $j$ collide? If so, when?

$$\Delta t = \begin{cases} \infty & \text{if } \Delta v \cdot \Delta r \geq 0 \\ \infty & \text{if } d < 0 \\ -\dfrac{\Delta v \cdot \Delta r + \sqrt{d}}{\Delta v \cdot \Delta v} & \text{otherwise} \end{cases}$$

$$d = (\Delta v \cdot \Delta r)^2 - (\Delta v \cdot \Delta v)(\Delta r \cdot \Delta r - \sigma^2) \qquad \sigma = \sigma_i + \sigma_j$$

$$\Delta v = (\Delta vx, \Delta vy) = (vx_i - vx_j, \ vy_i - vy_j)$$
$$\Delta r = (\Delta rx, \Delta ry) = (rx_i - rx_j, \ ry_i - ry_j)$$

$$\Delta v \cdot \Delta v = (\Delta vx)^2 + (\Delta vy)^2$$
$$\Delta r \cdot \Delta r = (\Delta rx)^2 + (\Delta ry)^2$$
$$\Delta v \cdot \Delta r = (\Delta vx)(\Delta rx) + (\Delta vy)(\Delta ry)$$

Important note: This is high-school physics, so we won't be testing you on it!

41

---

## Particle-particle collision resolution

Collision resolution.  When two particles collide, how does velocity change?

$$\begin{aligned} vx_i' &= vx_i + Jx / m_i \\ vy_i' &= vy_i + Jy / m_i \\ vx_j' &= vx_j - Jx / m_j \\ vy_j' &= vy_j - Jy / m_j \end{aligned}$$

Newton's second law
(momentum form)

$$Jx = \frac{J\,\Delta rx}{\sigma}, \quad Jy = \frac{J\,\Delta ry}{\sigma}, \quad J = \frac{2\,m_i\,m_j\,(\Delta v \cdot \Delta r)}{\sigma(m_i + m_j)}$$

impulse due to normal force
(conservation of energy, conservation of momentum)

Important note: This is high-school physics, so we won't be testing you on it!

42

---

## Particle data type skeleton

```
public class Particle
{
    private double rx, ry;        // position
    private double vx, vy;        // velocity
    private final double radius;  // radius
    private final double mass;    // mass
    private int count;            // number of collisions

    public Particle(...) { }

    public void move(double dt) { }
    public void draw()          { }

    public double timeToHit(Particle that)    { }
    public double timeToHitVerticalWall()     { }
    public double timeToHitHorizontalWall()   { }

    public void bounceOff(Particle that)      { }
    public void bounceOffVerticalWall()       { }
    public void bounceOffHorizontalWall()     { }

}
```

predict collision with
particle or wall

resolve collision with
particle or wall

43

---

## Particle-particle collision and resolution implementation

```
public double timeToHit(Particle that)
{
    if (this == that) return INFINITY;
    double dx  = that.rx - this.rx, dy  = that.ry - this.ry;
    double dvx = that.vx - this.vx; dvy = that.vy - this.vy;
    double dvdr = dx*dvx + dy*dvy;
    if( dvdr > 0) return INFINITY;
    double dvdv = dvx*dvx + dvy*dvy;
    double drdr = dx*dx + dy*dy;
    double sigma = this.radius + that.radius;
    double d = (dvdr*dvdr) - dvdv * (drdr - sigma*sigma);
    if (d < 0) return INFINITY;
    return -(dvdr + Math.sqrt(d)) / dvdv;
}
```

no collision

```
public void bounceOff(Particle that)
{
    double dx  = that.rx - this.rx, dy  = that.ry - this.ry;
    double dvx = that.vx - this.vx, dvy = that.vy - this.vy;
    double dvdr = dx*dvx + dy*dvy;
    double dist = this.radius + that.radius;
    double J = 2 * this.mass * that.mass * dvdr / ((this.mass + that.mass) * dist);
    double Jx = J * dx / dist;
    double Jy = J * dy / dist;
    this.vx += Jx / this.mass;
    this.vy += Jy / this.mass;
    that.vx -= Jx / that.mass;
    that.vy -= Jy / that.mass;
    this.count++;
    that.count++;
}
```

Important note: This is high-school physics, so we won't be testing you on it!

44

## Collision system: event-driven simulation main loop

### Initialization.
- Fill PQ with all potential particle-wall collisions.
- Fill PQ with all potential particle-particle collisions.

"potential" since collision may not happen if some other collision intervenes

two particles on a collision course

third particle interferes: no collision

An invalidated event

### Main loop.
- Delete the impending event from PQ (min priority = $t$).
- If the event has been invalidated, ignore it.
- Advance all particles to time $t$, on a straight-line trajectory.
- Update the velocities of the colliding particle(s).
- Predict future particle-wall and particle-particle collisions involving the colliding particle(s) and insert events onto PQ.

## Event data type

### Conventions.
- Neither particle `null` ⇒ particle-particle collision.
- One particle `null` ⇒ particle-wall collision.
- Both particles `null` ⇒ redraw event.

```
private class Event implements Comparable<Event>
{
    private double time;         // time of event
    private Particle a, b;       // particles involved in event
    private int countA, countB;  // collision counts for a and b

    public Event(double t, Particle a, Particle b) { }          ← create event

    public int compareTo(Event that)
    {    return this.time - that.time;    }                      ← ordered by time

    public boolean isValid()
    {    }                                                        ← invalid if intervening
                                                                    collision
}
```

## Collision system implementation: skeleton

```
public class CollisionSystem
{
    private MinPQ<Event> pq;       // the priority queue
    private double t  = 0.0;       // simulation clock time
    private Particle[] particles;  // the array of particles

    public CollisionSystem(Particle[] particles) { }

    private void predict(Particle a)
    {                                            add to PQ all particle-wall and particle-
        if (a == null) return;                   particle collisions involving this particle
        for (int i = 0; i < N; i++)
        {
            double dt = a.timeToHit(particles[i]);
            pq.insert(new Event(t + dt, a, particles[i]));
        }
        pq.insert(new Event(t + a.timeToHitVerticalWall()   , a, null));
        pq.insert(new Event(t + a.timeToHitHorizontalWall(), null, a));
    }

    private void redraw()  { }

    public void simulate() {  /* see next slide */  }
}
```

## Collision system implementation: main event-driven simulation loop

```
public void simulate()
{
    pq = new MinPQ<Event>();                                    ← initialize PQ with
    for(int i = 0; i < N; i++) predict(particles[i]);            collision events and
    pq.insert(new Event(0, null, null));                         redraw event

    while(!pq.isEmpty())
    {
        Event event = pq.delMin();
        if(!event.isValid()) continue;                          ← get next event
        Particle a = event.a;
        Particle b = event.b;

        for(int i = 0; i < N; i++)
            particles[i].move(event.time - t);                  ← update positions
        t = event.time;                                            and time

        if      (a != null && b != null) a.bounceOff(b);        ← process event
        else if (a != null && b == null) a.bounceOffVerticalWall()
        else if (a == null && b != null) b.bounceOffHorizontalWall();
        else if (a == null && b == null) redraw();

        predict(a);                                             ← predict new events
        predict(b);                                               based on changes
    }
}
```

## Simulation example 1

% java CollisionSystem 100

## Simulation example 2

% java CollisionSystem < billiards.txt

## Simulation example 3

% java CollisionSystem < brownian.txt

## Simulation example 4

% java CollisionSystem < diffusion.txt

‣ API

‣ sequential search

‣ binary search

‣ ordered operations

## Symbol tables

Key-value pair abstraction.

- Insert a value with specified key.
- Given a key, search for the corresponding value.

Ex. DNS lookup.

- Insert URL with specified IP address.
- Given URL, find corresponding IP address.

| URL | IP address |
|:---:|:---:|
| www.cs.princeton.edu | 128.112.136.11 |
| www.princeton.edu | 128.112.128.15 |
| www.yale.edu | 130.132.143.21 |
| www.harvard.edu | 128.103.060.55 |
| www.simpsons.com | 209.052.165.60 |

key ↑          value ↑

# Symbol table applications

| application | purpose of search | key | value |
| --- | --- | --- | --- |
| dictionary | find definition | word | definition |
| book index | find relevant pages | term | list of page numbers |
| file share | find song to download | name of song | computer ID |
| financial account | process transactions | account number | transaction details |
| web search | find relevant web pages | keyword | list of page names |
| compiler | find properties of variables | variable name | type and value |
| routing table | route Internet packets | destination | best route |
| DNS | find IP address given URL | URL | IP address |
| reverse DNS | find URL given IP address | IP address | URL |
| genomics | find markers | DNA string | known positions |
| file system | find file on disk | filename | location on disk |

# Symbol table API

Associative array abstraction.  Associate one value with each key.

```
public class ST<Key, Value>
```

|  |  |
|---|---|
| ST() | *create a symbol table* |
| void put(Key key, Value val) | *put key-value pair into the table (remove key from table if value is null)* |
| Value get(Key key) | *value paired with key (null if key is absent)* |
| void delete(Key key) | *remove key (and its value) from table* |
| boolean contains(Key key) | *is there a value paired with key?* |
| boolean isEmpty() | *is the table empty?* |
| int size() | *number of key-value pairs in the table* |
| Iterable<Key> keys() | *all the keys in the table* |

⟵ `a[key] = val;`

⟵ `a[key]`

**API for a generic basic symbol table**

## Conventions

- Values are not `null`.
- Method `get()` returns `null` if key not present.
- Method `put()` overwrites old value with new value.

**Intended consequences.**
- Easy to implement `contains()`.

```
public boolean contains(Key key)
{   return get(key) != null;   }
```

- Can implement lazy version of `delete()`.

```
public void delete(Key key)
{   put(key, null);             }
```

## Keys and values

**Value type.** Any generic type.

**Key type: several natural assumptions.**

- Assume keys are `Comparable`, use `compareTo()`.
- Assume keys are any generic type, use `equals()` to test equality.
- Assume keys are any generic type, use `equals()` to test equality and `hashCode()` to scramble key.

**Best practices.** Use immutable types for symbol table keys.

- Immutable in Java: `String, Integer, Double, File, …`
- Mutable in Java: `Date, StringBuilder, Url, …`

## ST test client for traces

Build ST by associating value i with ith string from standard input.

```
public static void main(String[] args)
{
    ST<String, Integer> st = new ST<String, Integer>();
    String[] a = StdIn.readAll().split("\\s+");
    for (int i = 0; i < a.length; i++)
        st.put(a[i], i);
    for (String s : st.keys())
        StdOut.println(s + " " + st.get(s));
}
```

output

keys    S  E  A  R  C  H  E  X  A  M  P  L  E

values  0  1  2  3  4  5  6  7  8  9 10 11 12

| A | 8 |
|---|---|
| C | 4 |
| E | 12 |
| H | 5 |
| L | 9 |
| M | 11 |
| P | 10 |
| R | 3 |
| S | 0 |
| X | 7 |

## ST test client for analysis

Frequency counter.  Read a sequence of strings from standard input and print out one that occurs with highest frequency.

```
% more tinyTale.txt
it was the best of times
it was the worst of times
it was the age of wisdom
it was the age of foolishness
it was the epoch of belief
it was the epoch of incredulity
it was the season of light
it was the season of darkness
it was the spring of hope
it was the winter of despair

% java FrequencyCounter 1 < tinyTale.txt          ⟵——— tiny example (60 words, 20 distinct)
it 10

% java FrequencyCounter 8 < tale.txt              ⟵——— real example (135,635 words, 10,769 distinct)
business 122

% java FrequencyCounter 10 < leipzig1M.txt  ⟵——— real example (21,191,455 words, 534,580 distinct)
government 24763
```

## Frequency counter implementation

```java
public class FrequencyCounter
{
   public static void main(String[] args)
   {
      int minlen = Integer.parseInt(args[0]);
      ST<String, Integer> st = new ST<String, Integer>();        create ST
      while (!StdIn.isEmpty())
      {
         String word = StdIn.readString();          ignore short strings
         if (word.length() < minlen) continue;       read string and
         if (!st.contains(word)) st.put(word, 1);     update frequency
         else                    st.put(word, st.get(word) + 1);
      }
      String max = "";
      st.put(max, 0);
      for (String word : st.keys())                   print a string
         if (st.get(word) > st.get(max))              with max freq
            max = word;
      StdOut.println(max + " " + st.get(max));
   }
}
```

9

## Sequential search in a linked list

Data structure. Maintain an (unordered) linked list of key-value pairs.

Search. Scan through all keys until find a match.

Insert. Scan through all keys until find a match; if no match add to front.



**Trace of linked-list ST implementation for standard indexing client**

# Elementary ST implementations:  summary

| ST implementation | worst case | | average case | | ordered iteration? | operations on keys |
|---|---|---|---|---|---|---|
| | search | insert | search hit | insert | | |
| sequential search (unordered list) | $N$ | $N$ | $N / 2$ | $N$ | no | `equals()` |



**Costs for** `java FrequencyCounter 8 < tale.txt` **using** `LinkedListST`

**Challenge.**  Efficient implementations of both search and insert.

# Binary search

**Data structure.** Maintain an ordered array of key-value pairs.

**Rank helper function.** How many keys < k?

```
                                keys[]

successful search for P     0   1   2   3   4   5   6   7   8   9

        lo  hi  m
        0   9   4       A   C   E   H   L   M   P   R   S   X          entries in black
                                                                       are a[lo..hi]
        5   9   7       A   C   E   H   L   M   P   R   S   X

        5   6   5       A   C   E   H   L   M   P   R   S   X
                                                                   entry in red is a[m]
        6   6   6       A   C   E   H   L   M   P   R   S   X

unsuccessful search for Q                          loop exits with keys[m] = P:  return 6

        lo  hi  m
        0   9   4       A   C   E   H   L   M   P   R   S   X

        5   9   7       A   C   E   H   L   M   P   R   S   X

        5   6   5       A   C   E   H   L   M   P   R   S   X

        7   6   6       A   C   E   H   L   M   P   R   S   X
       ___
            loop exits with lo > hi:  return  7
```

**Trace of binary search for rank in an ordered array**

## Binary search:  Java implementation

```
public Value get(Key key)
{
   if (isEmpty()) return null;
   int i = rank(key);
   if (i < N && keys[i].compareTo(key) == 0) return vals[i];
   else  return null;
}
```

```
private int rank(Key key)                      number of keys < key
{
   int lo = 0, hi = N-1;
   while (lo <= hi)
   {
       int mid = lo + (hi - lo) / 2;
       int cmp = key.compareTo(keys[mid]);
       if      (cmp  < 0) hi = mid - 1;
       else if (cmp  > 0) lo = mid + 1;
       else if (cmp == 0) return mid;
   }
   return lo;
}
```

Binary search: mathematical analysis

Proposition. Binary search uses $\sim \lg N$ compares to search any array of size $N$.

Def. $T(N) \equiv$ number of compares to binary search in a sorted array of size $N$.

$$\leq\ T(N/2)\ +\ 1$$

↑
left or right  half

Binary search recurrence. $T(N) \leq T(N/2) + 1$ for $N > 1$, with $T(1) = 1$.

- Not quite right for odd $N$.
- Same recurrence holds for many algorithms.

Solution. $T(N) \sim \lg N$.

- For simplicity, we'll prove when $N$ is a power of 2.
- True for all $N$. [see COS 340]

## Binary search recurrence

Binary search recurrence. $T(N) \le T(N / 2) + 1$ for $N > 1$, with $T(1) = 1$.

Proposition. If $N$ is a power of 2, then $T(N) \le \lg N + 1$.

Pf.

| | | |
|---|---|---|
| $T(N)$ | $\le T(N / 2) + 1$ | given |
| | $\le T(N / 4) + 1 + 1$ | apply recurrence to first term |
| | $\le T(N / 8) + 1 + 1 + 1$ | apply recurrence to first term |
| | $\ldots$ | |
| | $\le T(N / N) + 1 + 1 + \ldots + 1$ | stop applying, T(1) = 1 |
| | $= \lg N + 1$ | |

# Binary search: trace of standard indexing client

Problem. To insert, need to shift all greater keys over.



|  |  | keys[] |  |  |  |  |  |  |  |  |  | N | vals[] |  |  |  |  |  |  |  |  |  |  |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| key | value | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |  | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| S | 0 | S |  |  |  |  |  |  |  |  |  | 1 | 0 |  |  |  |  |  |  |  |  |  |  |
| E | 1 | E | S |  |  |  |  |  |  |  |  | 2 | 1 | 0 |  |  |  |  |  |  |  |  |  |
| A | 2 | A | E | S |  |  |  |  |  |  |  | 3 | 2 | 1 | 0 |  |  |  |  |  |  |  |  |
| R | 3 | A | E | R | S |  |  |  |  |  |  | 4 | 2 | 1 | 3 | 0 |  |  |  |  |  |  |  |
| C | 4 | A | C | E | R | S |  |  |  |  |  | 5 | 2 | 4 | 1 | 3 | 0 |  |  |  |  |  |  |
| H | 5 | A | C | E | H | R | S |  |  |  |  | 6 | 2 | 4 | 1 | 5 | 3 | 0 |  |  |  |  |  |
| E | 6 | A | C | E | H | R | S |  |  |  |  | 6 | 2 | 4 | (6) | 5 | 3 | 0 |  |  |  |  |  |
| X | 7 | A | C | E | H | R | S | X |  |  |  | 7 | 2 | 4 | 6 | 5 | 3 | 0 | 7 |  |  |  |  |
| A | 8 | A | C | E | H | R | S | X |  |  |  | 7 | (8) | 4 | 6 | 5 | 3 | 0 | 7 |  |  |  |  |
| M | 9 | A | C | E | H | M | R | S | X |  |  | 8 | 8 | 4 | 6 | 5 | 9 | 3 | 0 | 7 |  |  |  |
| P | 10 | A | C | E | H | M | P | R | S | X |  | 9 | 8 | 4 | 6 | 5 | 9 | 10 | 3 | 0 | 7 |  |  |
| L | 11 | A | C | E | H | L | M | P | R | S | X | 10 | 8 | 4 | 6 | 5 | 11 | 9 | 10 | 3 | 0 | 7 |
| E | 12 | A | C | E | H | L | M | P | R | S | X | 10 | 8 | 4 | (12) | 5 | 11 | 9 | 10 | 3 | 0 | 7 |
|  |  | A | C | E | H | L | M | P | R | S | X |  | 8 | 4 | 12 | 5 | 11 | 9 | 10 | 3 | 0 | 7 |

*entries in red were inserted*

*entries in gray did not move*

*entries in black moved to the right*

*circled entries are changed values*

# Elementary ST implementations:  summary

| ST implementation | worst case | | average case | | ordered iteration? | operations on keys |
|---|---|---|---|---|---|---|
| | search | insert | search hit | insert | | |
| sequential search (unordered list) | N | N | N / 2 | N | no | `equals()` |
| binary search (ordered array) | log N | N | log N | N / 2 | yes | `compareTo()` |



**Costs for** `java FrequencyCounter 8 < tale.txt` **using** `OrderedArrayST`

**Challenge.** Efficient implementations of both search and insert.

## Ordered symbol table API



**Examples of ordered symbol-table operations**

## Ordered symbol table API

```
public class ST<Key extends Comparable<Key>, Value>
```

|  |  |
|---|---|
| `ST()` | *create an ordered symbol table* |
| `void put(Key key, Value val)` | *put key-value pair into the table (remove* key *from table if value is null)* |
| `Value get(Key key)` | *value paired with* key *(null if* key *is absent)* |
| `void delete(Key key)` | *remove* key *(and its value) from table* |
| `boolean contains(Key key)` | *is there a value paired with* key*?* |
| `boolean isEmpty()` | *is the table empty?* |
| `int size()` | *number of key-value pairs* |
| `Key min()` | *smallest key* |
| `Key max()` | *largest key* |
| `Key floor(Key key)` | *largest key less than or equal to* key |
| `Key ceiling(Key key)` | *smallest key greater than or equal to* key |
| `int rank(Key key)` | *number of keys less than* key |
| `Key select(int k)` | *key of rank* k |
| `void deleteMin()` | *delete smallest key* |
| `void deleteMax()` | *delete largest key* |
| `int size(Key lo, Key hi)` | *number of keys in* `[lo..hi]` |
| `Iterable<Key> keys(Key lo, Key hi)` | *keys in* `[lo..hi], in sorted order` |
| `Iterable<Key> keys()` | *all keys in the table, in sorted order* |

**API for a generic ordered symbol table**

# Binary search:  ordered symbol table operations summary

| | sequential search | binary search |
|---|---|---|
| search | N | lg N |
| insert | 1 | N |
| min / max | N | 1 |
| floor / ceiling | N | lg N |
| rank | N | lg N |
| select | N | 1 |
| ordered iteration | N log N | N |

worst-case running time of ordered symbol table operations

# 3.2  Binary Search Trees



▸ BSTs

▸ ordered operations

▸ deletion

# Binary search trees

Definition. A BST is a binary tree in symmetric order.

A binary tree is either:
- Empty.
- Two disjoint binary trees (left and right).



**Anatomy of a binary tree**

Symmetric order.

Each node has a key, and every node's key is:
- Larger than all keys in its left subtree.
- Smaller than all keys in its right subtree.



**Anatomy of a binary search tree**

# BST representation in Java

Java definition.  A BST is a reference to a root Node.

A Node is comprised of four fields:

- A Key and a Value.
- A reference to the left and right subtree.

smaller keys          larger keys

```
private class Node
{
    private Key key;
    private Value val;
    private Node left, right;
    public Node(Key key, Value val)
    {
        this.key = key;
        this.val = val;
    }
}
```

Key and Value are generic types; Key is Comparable



Binary search tree

## BST implementation (skeleton)

```
public class BST<Key extends Comparable<Key>, Value>
{
    private Node root;                              ← root of BST

    private class Node
    {  /* see previous slide */  }

    public void put(Key key, Value val)
    {  /* see next slides */  }

    public Value get(Key key)
    {  /* see next slides */  }

    public void delete(Key key)
    {  /* see next slides */  }

    public Iterable<Key> iterator()
    {  /* see next slides */  }

}
```

# BST search

Get. Return value corresponding to given key, or `null` if no such key.



**successful search for** R

R *is less than S
so look to the left*

*black nodes could
match the search key*

R *is greater than E
so look to the right*

*gray nodes cannot
match the search key*

*found R
(search hit)
so return value*

**unsuccessful search for** T

T *is greater than S
so look to the right*

T *is less than X
so look to the left*

*link is null
so T is not in tree
(search miss)*

## BST search: Java implementation

Get. Return value corresponding to given key, or `null` if no such key.

```java
public Value get(Key key)
{
   Node x = root;
   while (x != null)
   {
      int cmp = key.compareTo(x.key);
      if        (cmp  < 0) x = x.left;
      else if (cmp  > 0) x = x.right;
      else if (cmp == 0) return x.val;
   }
   return null;
}
```

Running time. Proportional to depth of node.

# BST insert

Put.  Associate value with key.

Search for key, then two cases:
- Key in tree  $\Rightarrow$  reset value.
- Key not in tree $\Rightarrow$  add new node.



inserting L

search for L ends
at this null link

create new node

reset links
on the way up

**Insertion into a BST**

# BST insert: Java implementation

Put.  Associate value with key.

```
public void put(Key key, Value val)
{   root = put(root, key, val);   }


private Node put(Node x, Key key, Value val)
{
    if (x == null) return new Node(key, val);
    int cmp = key.compareTo(x.key);
    if        (cmp  < 0)
        x.left  = put(x.left,  key, val);
    else if (cmp  > 0)
        x.right = put(x.right, key, val);
    else if (cmp == 0)
        x.val = val;

    return x;

}
```

concise, but tricky,
recursive code;
read carefully!

Running time.  Proportional to depth of node.

# BST trace: standard indexing client

## Tree shape

- Many BSTs correspond to same set of keys.
- Cost of search/insert is proportional to depth of node.



**Remark.** Tree shape depends on order of insertion.

## BST insertion: random order

Observation. If keys inserted in random order, tree stays relatively flat.



N = 255
max = 16
avg = 9.1
opt = 7.0

# BST insertion: random order visualization

Ex.  Insert keys in random order.



N = 255
max = 16
avg = 9.1
opt = 7.0

# Correspondence between BSTs and quicksort partitioning



**Remark.** Correspondence is 1-1 if no duplicate keys.

BSTs: mathematical analysis

Proposition. If keys are inserted in random order, the expected number of compares for a search/insert is ~ 2 ln N.

Pf. 1-1 correspondence with quicksort partitioning.

Proposition. [Reed, 2003] If keys are inserted in random order, expected height of tree is ~ 4.311 ln N.

But... Worst-case for search/insert/height is N.
(exponentially small chance when keys are inserted in random order)

## ST implementations: summary

| implementation | guarantee | | average case | | ordered ops? | operations on keys |
| --- | --- | --- | --- | --- | --- | --- |
| | search | insert | search hit | insert | | |
| sequential search (unordered list) | N | N | N/2 | N | no | `equals()` |
| binary search (ordered array) | lg N | N | lg N | N/2 | yes | `compareTo()` |
| BST | N | N | 1.39 lg N | 1.39 lg N | ? | `compareTo()` |



**Costs for** java FrequencyCounter 8 < tale.txt **using** BST

▶ BSTs

▶ **ordered operations**

▶ deletion

# Minimum and maximum

**Minimum.** Smallest key in table.

**Maximum.** Largest key in table.



**Q.** How to find the min / max.

# Floor and ceiling

**Floor.** Largest key ≤ to a given key.

**Ceiling.** Smallest key ≥ to a given key.



**Q.** How to find the floor /ceiling.

# Computing the floor

**Case 1.** [k equals the key at root]

The floor of k is k.


**Case 2.** [k is less than the key at root]

The floor of k is in the left subtree.


**Case 3.** [k is greater than the key at root]

The floor of k is in the right subtree
(if there is any key ≤ k in right subtree);
otherwise it is the key in the root.



finding `floor(G)`

*G is less than S so*
`floor(G)` *must be
on the left*

*G is greater than E so*
`floor(G)` *could be
on the right*

`floor(G)` *in left
subtree is* `null`

*result*

19

# Computing the floor

```
public Key floor(Key key)
{
    Node x = floor(root, key);
    if (x == null) return null;
    return x.key;
}
private Node floor(Node x, Key key)
{
    if (x == null) return null;
    int cmp = key.compareTo(x.key);

    if (cmp == 0) return x;

    if (cmp < 0)  return floor(x.left, key);

    Node t = floor(x.right, key);
    if (t != null) return t;
    else           return x;

}
```



finding `floor(G)`

G *is less than S so* `floor(G)` *must be on the left*

G *is greater than E so* `floor(G)` *could be on the right*

`floor(G)` *in left subtree is* `null`

*result*

## Subtree counts

In each node, we store the number of nodes in the subtree rooted at that node.
To implement `size()`, return the count at the root.



**Remark.** This facilitates efficient implementation of `rank()` and `select()`.

## BST implementation:  subtree counts

```
private class Node
{
    private Key key;
    private Value val;
    private Node left;
    private Node right;
    private int N;
}
```

nodes in subtree

```
public int size()
{   return size(root);   }
```

```
private int size(Node x)
{
    if (x == null) return 0;
    return x.N;
}
```

```
private Node put(Node x, Key key, Value val)
{
    if (x == null) return new Node(key, val);
    int cmp = key.compareTo(x.key);
    if      (cmp  < 0) x.left  = put(x.left,  key, val);
    else if (cmp  > 0) x.right = put(x.right, key, val);
    else if (cmp == 0) x.val = val;
    x.N = 1 + size(x.left) + size(x.right);
    return x;
}
```

# Rank

Rank.  How many keys < k?

Easy recursive algorithm (4 cases!)



node count N

```
public int rank(Key key)
{   return rank(key, root);   }


private int rank(Key key, Node x)
{
    if (x == null) return 0;
    int cmp = key.compareTo(x.key);
    if      (cmp < 0) return rank(key, x.left);
    else if (cmp > 0) return 1 + size(x.left) + rank(key, x.right);
    else              return size(x.left);
}
```

## Inorder traversal

- Traverse left subtree.
- Enqueue key.
- Traverse right subtree.

```java
public Iterable<Key> keys()
{
    Queue<Key> q = new Queue<Key>();
    inorder(root, queue);
    return q;
}


private void inorder(Node x, Queue<Key> q)
{
    if (x == null) return;
    inorder(x.left, q);
    q.enqueue(x.key);
    inorder(x.right, q);
}
```



BST

key  val

left    right

BST with smaller keys          BST with larger keys

smaller keys, in order    key    larger keys, in order

all keys, in order

**Property.** Inorder traversal of a BST yields keys in ascending order.

# Inorder traversal

- Traverse left subtree.
- Enqueue key.
- Traverse right subtree.



```
inorder(S)
   inorder(E)
      inorder(A)
         enqueue A
         inorder(C)
            enqueue C
      enqueue E
      inorder(R)
         inorder(H)
            enqueue H
            inorder(M)
               enqueue M
         print R
   enqueue S
   inorder(X)
      enqueue X
```

recursive calls

```
A

C
E

H

M
R
S

X
```

queue

```
S
S E
S E A

S E A C



S E R
S E R H

S E R H M




S X
```

function call stack

A C E H M R S X

# BST: ordered symbol table operations summary

| | sequential search | binary search | BST |
|---|---|---|---|
| search | N | lg N | h |
| insert | 1 | N | h |
| min / max | N | 1 | h |
| floor / ceiling | N | lg N | h |
| rank | N | lg N | h |
| select | N | 1 | h |
| ordered iteration | N log N | N | N |

h = height of BST
(proportional to log N
if keys inserted in random order)

worst-case running time of ordered symbol table operations

▸ BSTs

▸ ordered operations

▸ **deletion**

## ST implementations: summary

| implementation | guarantee | | | average case | | | ordered iteration? | operations on keys |
|---|---|---|---|---|---|---|---|---|
| | search | insert | delete | search hit | insert | delete | | |
| sequential search (linked list) | N | N | N | N/2 | N | N/2 | no | `equals()` |
| binary search (ordered array) | lg N | N | N | lg N | N/2 | N/2 | yes | `compareTo()` |
| BST | N | N | N | 1.39 lg N | 1.39 lg N | ??? | yes | `compareTo()` |

Next.  Deletion in BSTs.

## BST deletion: lazy approach

To remove a node with a given key:

- Set its value to `null`.

- Leave key in tree to guide searches (but don't consider it equal to search key).



**Cost.** O(log N') per insert, search, and delete (if keys in random order), where N' is the number of key-value pairs ever inserted in the BST.

**Unsatisfactory solution.** Tombstone overload.

## Deleting the minimum

To delete the minimum key:

- Go left until finding a node with a null left link.
- Replace that node by its right link.
- Update subtree counts.

```
public void deleteMin()
{   root = deleteMin(root);   }


private Node deleteMin(Node x)
{
    if (x.left == null) return x.right;
    x.left = deleteMin(x.left);
    x.N = 1 + size(x.left) + size(x.right);
    return x;
}
```



*go left until reaching null left link*

*return that node's right link*

*available for garbage collection*

*update links and counts after recursive calls*

# Hibbard deletion

To delete a node with key k: search for node t containing key k.

Case 0. [0 children] Delete t by setting parent link to null.



deleting C

node to delete

replace with
null link

available for
garbage
collection

update counts after
recursive calls

# Hibbard deletion

To delete a node with key k:  search for node t containing key k.

Case 1.  [1 child]  Delete t by replacing parent link.



deleting R

node to delete

replace with child link

available for garbage collection

update counts after recursive calls

# Hibbard deletion

To delete a node with key k:  search for node t containing key k.

## Case 2.  [2 children]

- Find successor x of t.          ⟵ x has no left child
- Delete the minimum in t's right subtree.    ⟵ but don't garbage collect x
- Put x in t's spot.          ⟵ still a BST



deleting E

node to delete

E

search for key E

x

t.left    deleteMin(t.right)

t

x    successor
min(t.right)

go right, then
go left until
reaching null
left link

update links and
node counts after
recursive calls

# Hibbard deletion:  Java implementation

```java
public void delete(Key key)
{   root = delete(root, key);   }

private Node delete(Node x, Key key) {
    if (x == null) return null;
    int cmp = key.compareTo(x.key);
    if      (cmp < 0) x.left  = delete(x.left,  key);    ← search for key
    else if (cmp > 0) x.right = delete(x.right, key);
    else {
        if (x.right == null) return x.left;              ← no right child

        Node t = x;
        x = min(t.right);
        x.right = deleteMin(t.right);                    ← replace with successor
        x.left = t.left;
    }
    x.N = size(x.left) + size(x.right) + 1;              ← update subtree counts
    return x;
}
```

## Hibbard deletion: analysis

**Unsatisfactory solution.** Not symmetric.



N = 150
max = 16
avg = 9.3
opt = 6.4

**Surprising consequence.** Trees not random (!) ⇒ sqrt(N) per op.

**Longstanding open problem.** Simple and efficient delete for BSTs.

# ST implementations:  summary

| implementation | guarantee | | | average case | | | ordered iteration? | operations on keys |
|---|---|---|---|---|---|---|---|---|
| | search | insert | delete | search hit | insert | delete | | |
| sequential search (linked list) | N | N | N | N/2 | N | N/2 | no | `equals()` |
| binary search (ordered array) | lg N | N | N | lg N | N/2 | N/2 | yes | `compareTo()` |
| BST | N | N | N | 1.39 lg N | 1.39 lg N | √N | yes | `compareTo()` |

other operations also become √N
if deletions allowed

Next lecture.  Guarantee logarithmic performance for all operations.

# 3.3 Balanced Trees



- ‣ **2-3 trees**
- ‣ **red-black trees**
- ‣ **B-trees**

# Symbol table review

| implementation | guarantee | | | average case | | | ordered iteration? | operations on keys |
|---|---|---|---|---|---|---|---|---|
| | search | insert | delete | search hit | insert | delete | | |
| sequential search (linked list) | N | N | N | N/2 | N | N/2 | no | `equals()` |
| binary search (ordered array) | lg N | N | N | lg N | N/2 | N/2 | yes | `compareTo()` |
| BST | N | N | N | 1.39 lg N | 1.39 lg N | ? | yes | `compareTo()` |
| Goal | log N | log N | log N | log N | log N | log N | yes | `compareTo()` |

Challenge.  Guarantee performance.

This lecture.  2-3 trees, left-leaning red-black trees, B-trees.

introduced to the world in
COS 226, Fall 2007

‣ **2-3 trees**

‣ red-black trees

‣ B-trees

## 2-3 tree

Allow 1 or 2 keys per node.

- 2-node: one key, two children.
- 3-node: two keys, three children.

Symmetric order. Inorder traversal yields keys in ascending order.

Perfect balance. Every path from root to null link has same length.

# Search in a 2-3 tree

- Compare search key against keys in node.
- Find interval containing search key.
- Follow associated link (recursively).



**successful search for** H

H *is less than M so look to the left* → M

H *is between E and L so look in the middle*

*found H so return value (search hit)*

**unsuccessful search for** B

B *is less than M so look to the left* → M

B *is less than E so look to the left*

B *is between A and C so look in the middle link is null so B is not in the tree (search miss)*

# Insertion in a 2-3 tree

**Case 1.** Insert into a 2-node at bottom.

- Search for key, as usual.
- Replace 2-node with 3-node.



inserting K

*search for K ends here*

*replace 2-node with new 3-node containing K*

# Insertion in a 2-3 tree

**Case 2.** Insert into a 3-node at bottom.

- Add new key to 3-node to create temporary 4-node.
- Move middle key in 4-node into parent.

why middle key?



inserting Z

search for Z ends at this 3-node

replace 3-node with temporary 4-node containing Z

replace 2-node with new 3-node containing middle key

split 4-node into two 2-nodes pass middle key to parent

# Insertion in a 2-3 tree

Case 2. Insert into a 3-node at bottom.

- Add new key to 3-node to create temporary 4-node.
- Move middle key in 4-node into parent.
- Repeat up the tree, as necessary.



inserting D

search for D ends
at this 3-node

add new key D to 3-node
to make temporary 4-node

add middle key C to 3-node
to make temporary 4-node

split 4-node into two 2-nodes
pass middle key to parent

add middle key E to 2-node
to make new 3-node

split 4-node into two 2-nodes
pass middle key to parent

# Insertion in a 2-3 tree

**Case 2.** Insert into a 3-node at bottom.

- Add new key to 3-node to create temporary 4-node.
- Move middle key in 4-node into parent.
- Repeat up the tree, as necessary.
- If you reach the root and it's a 4-node, split it into three 2-nodes.



**Remark.** Splitting the root increases height by 1.

# 2-3 tree construction trace

Standard indexing client.

# 2-3 tree construction trace

The same keys inserted in ascending order.

# Local transformations in a 2-3 tree

Splitting a 4-node is a local transformation:  constant number of operations.

# Global properties in a 2-3 tree

**Invariant.** Symmetric order.

**Invariant.** Perfect balance.

**Pf.** Each transformation maintains order and balance.

## 2-3 tree:  performance

Perfect balance.  Every path from root to null link has same length.



Tree height.

- Worst case:
- Best case:

## 2-3 tree: performance

**Perfect balance.** Every path from root to null link has same length.



**Tree height.**

- Worst case:    lg N.                                    [all 2-nodes]
- Best case:       $\log_3 N \approx .631$ lg N.    [all 3-nodes]
- Between 12 and 20 for a million nodes.
- Between 18 and 30 for a billion nodes.

Guaranteed logarithmic performance for search and insert.

# ST implementations:  summary

| implementation | guarantee | | | average case | | | ordered iteration? | operations on keys |
|---|---|---|---|---|---|---|---|---|
| | search | insert | delete | search hit | insert | delete | | |
| sequential search (linked list) | N | N | N | N/2 | N | N/2 | no | `equals()` |
| binary search (ordered array) | lg N | N | N | lg N | N/2 | N/2 | yes | `compareTo()` |
| BST | N | N | N | 1.39 lg N | 1.39 lg N | ? | yes | `compareTo()` |
| 2-3 tree | c lg N | c lg N | c lg N | c lg N | c lg N | c lg N | yes | `compareTo()` |

constants depend upon
implementation

## 2-3 tree: implementation?

Direct implementation is complicated, because:

- Maintaining multiple node types is cumbersome.
- Need multiple compares to move down tree.
- Need to move back up the tree to split 4-nodes.
- Large number of cases for splitting.

Bottom line.  Could do it, but there's a better way.

- 2-3-4 trees
- **red-black trees**
- B-trees

# Left-leaning red-black trees (Guibas-Sedgewick 1979 and Sedgewick 2007)

1. Represent 2–3 tree as a BST.

2. Use "internal" left-leaning links as "glue" for 3–nodes.



**3-node**: a b — less than a, between a and b, greater than b

larger key is root — b, a: less than a, between a and b, greater than b



2-3 tree

red links "glue" nodes within a 3-node

black links connect 2-nodes and 3-nodes

red-black tree

## An equivalent definition

A BST such that:

- No node has two red links connected to it.
- Every path from root to null link has the same number of black links.
- Red links lean left.

"perfect black balance"

# Left-leaning red-black trees: 1-1 correspondence with 2-3 trees

**Key property.** 1–1 correspondence between 2–3 and LLRB.



red−black tree

horizontal red links

2-3 tree

# Search implementation for red-black trees

Observation. Search is the same as for elementary BST (ignore color).

↑

but runs faster because of better balance

```java
public Val get(Key key)
{
    Node x = root;
    while (x != null)
    {
        int cmp = key.compareTo(x.key);
        if       (cmp  < 0) x = x.left;
        else if (cmp  > 0) x = x.right;
        else if (cmp == 0) return x.val;
    }
    return null;
}
```



Remark. Many other ops (e.g., ceiling, selection, iteration) are also identical.

# Red-black tree representation

Each node is pointed to by precisely one link (from its parent) ⇒
can encode color of links in nodes.

```
private static final boolean RED   = true;
private static final boolean BLACK = false;

private class Node
{
    Key key;
    Value val;
    Node left, right;
    boolean color;   // color of parent link
}

private boolean isRed(Node x)
{
    if (x == null) return false;
    return x.color == RED;
}
```

null links are black

h.left.color *is* RED

h

h.right.color *is* BLACK

# Elementary red-black tree operations

**Left rotation.** Orient a (temporarily) right-leaning red link to lean left.



```
private Node rotateLeft(Node h)
{
    assert (h != null) && isRed(h.right);
    Node x = h.right;
    h.right = x.left;
    x.left = h;
    x.color = h.color;
    h.color = RED;
    return x;
}
```

**Invariants.** Maintains symmetric order and perfect black balance.

# Elementary red-black tree operations

Right rotation.   Orient a left-leaning red link to (temporarily) lean right.



```
private Node rotateRight(Node h)
{
    assert (h != null) && isRed(h.left);
    Node x = h.left;
    h.left = x.right;
    x.right = h;
    x.color = h.color;
    h.color = RED;
    return x;
}
```

Invariants.  Maintains symmetric order and perfect black balance.

# Elementary red-black tree operations

**Color flip.** Recolor to split a (temporary) 4-node.



```
private void flipColors(Node h)
{
    assert !isRed(h) && isRed(h.left) && isRed(h.right);

    h.color = RED;
    h.left.color = BLACK;
    h.right.color = BLACK;
}
```

**Invariants.** Maintains symmetric order and perfect black balance.

# Insertion in a LLRB tree:  overview

Basic strategy.  Maintain 1-1 correspondence with 2-3 trees by applying elementary red-black tree operations

# Insertion in a LLRB tree

Warmup 1. Insert into a tree with exactly 1 node.



**left**

root

b

*search ends
at this null link*

root

b

a

*red link to
new node
containing* a
*converts 2-node
to 3-node*

**right**

root

a

*search ends
at this null link*

a

b

*attached new node
with red link*

root

b

a

*rotated left
to make a
legal 3-node*

# Insertion in a LLRB tree

Case 1.  Insert into a 2-node at the bottom.

• Do standard BST insert; color new link red.

• If new red link is a right link, rotate left.

# Insertion in a LLRB tree

**Warmup 2.** Insert into a tree with exactly 2 nodes.

# Insertion in a LLRB tree

**Case 2.** Insert into a 3-node at the bottom.

- Do standard BST insert; color new link red.
- Rotate to balance the 4-node (if needed).
- Flip colors to pass red link up one level.
- Rotate to make lean left (if needed).

# Insertion in a LLRB tree: passing red links up the tree

Case 2. Insert into a 3-node at the bottom.

- Do standard BST insert; color new link red.
- Rotate to balance the 4-node (if needed).
- Flip colors to pass red link up one level.
- Rotate to make lean left (if needed).
- Repeat Case 1 or Case 2 up the tree (if needed).

# LLRB tree construction trace

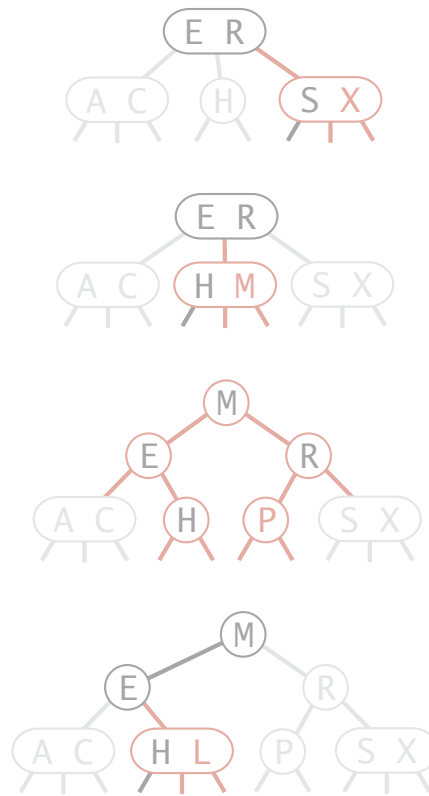Standard indexing client.



red black tree            2-3 tree

# LLRB tree construction trace

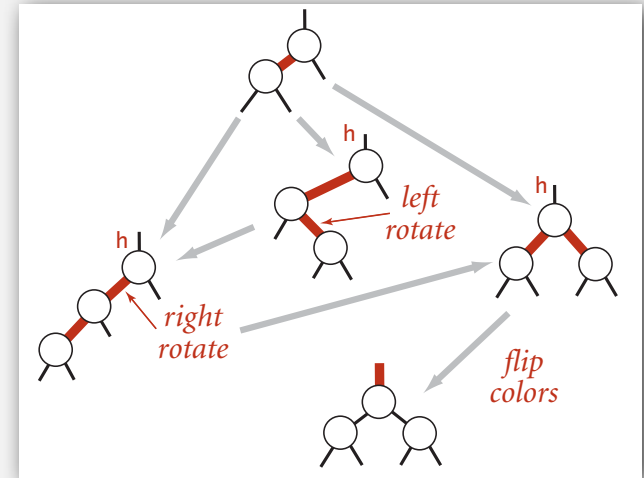Standard indexing client (continued).



red black tree                                    2-3 tree

# Insertion in a LLRB tree: Java implementation

Same code for both cases.

- Right child red, left child black: rotate left.
- Left child, left-left grandchild red: rotate right.
- Both children red: flip colors.



```java
private Node put(Node h, Key key, Value val)
{
    if (h == null) return new Node(key, val, RED);         ← insert at bottom
    int cmp = key.compareTo(h.key);
    if      (cmp < 0) h.left  = put(h.left,  key, val);
    else if (cmp > 0) h.right = put(h.right, key, val);
    else h.val = val;

    if (isRed(h.right) && !isRed(h.left))     h = rotateLeft(h);    ← lean left
    if (isRed(h.left)  && isRed(h.left.left)) h = rotateRight(h);   ← balance 4-node
    if (isRed(h.left)  && isRed(h.right))     h = flipColors(h);    ← split 4-node

    return h;
}
```
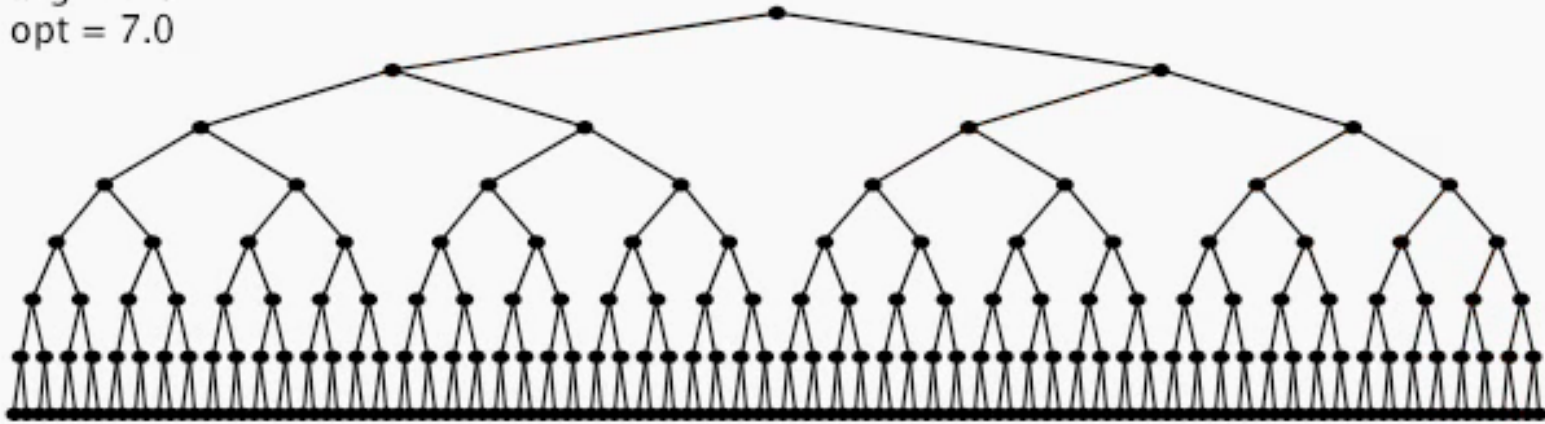
only a few extra lines of code
to provide near-perfect balance

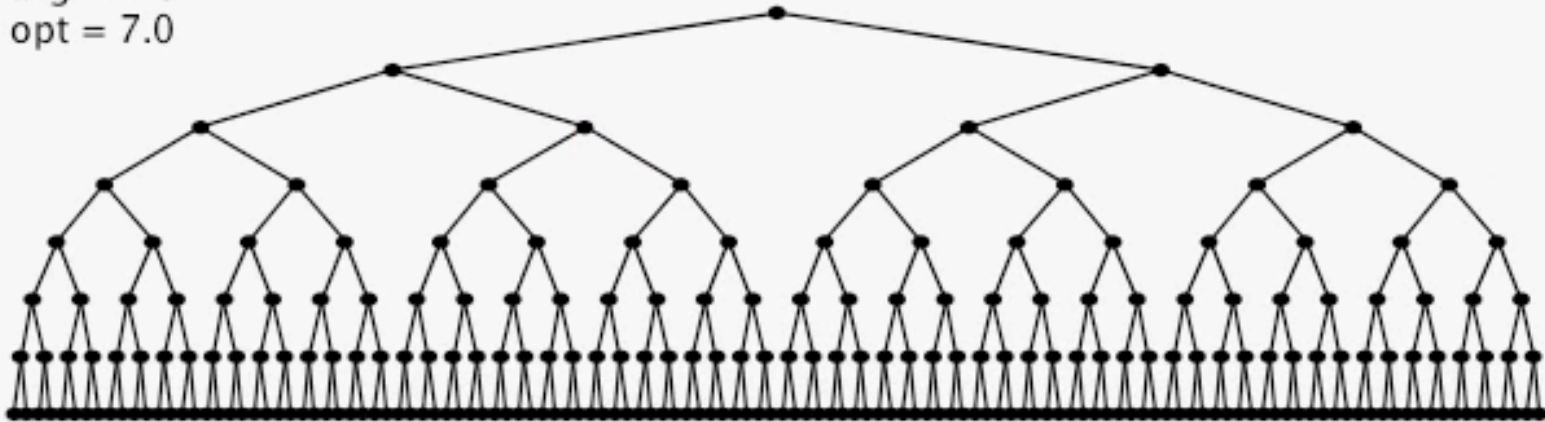# Insertion in a LLRB tree:  visualization

N = 255
max = 8
avg = 7.0
opt = 7.0

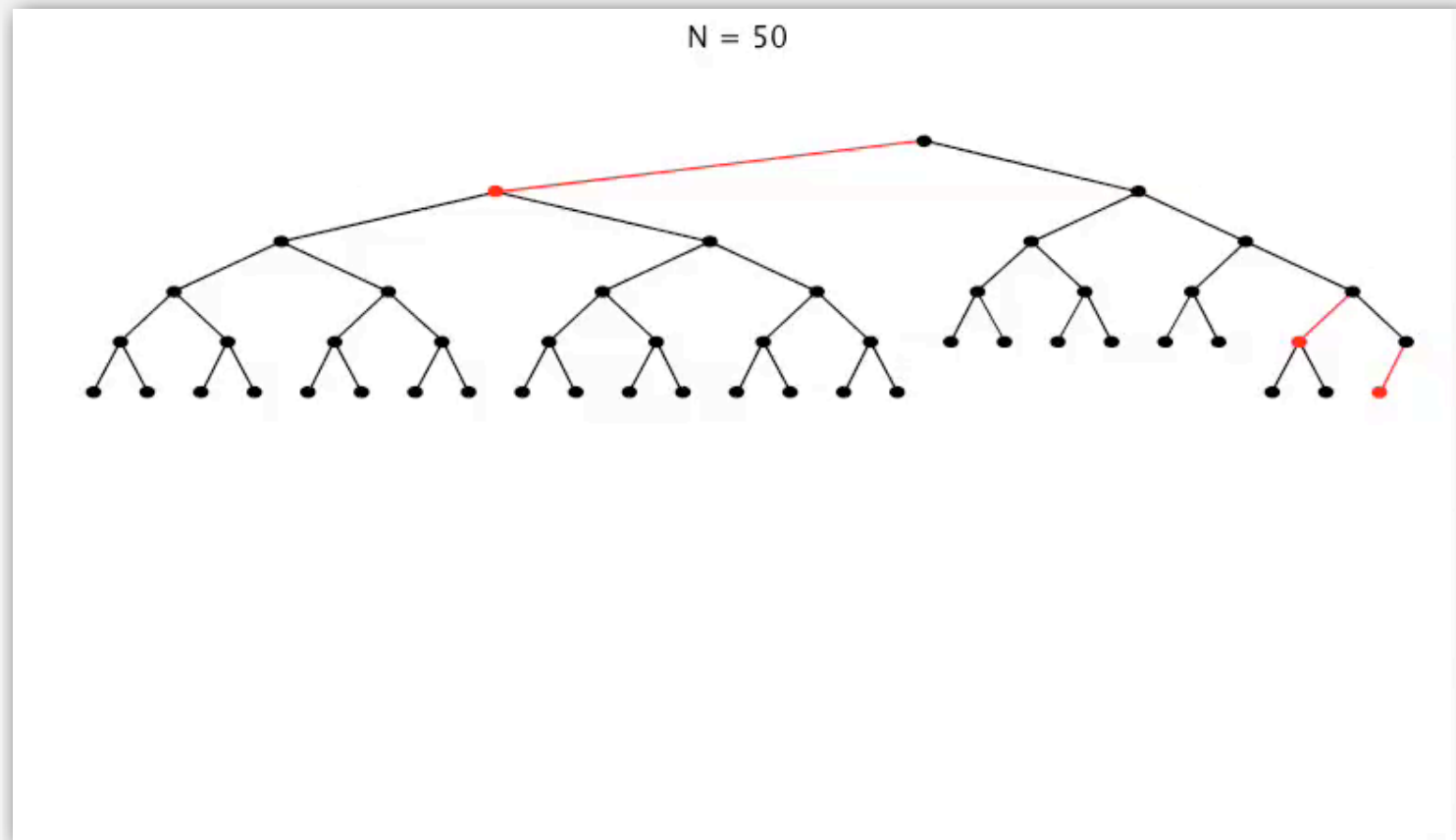255 insertions in ascending order

# Insertion in a LLRB tree: visualization



N = 255
max = 8
avg = 7.0
opt = 7.0
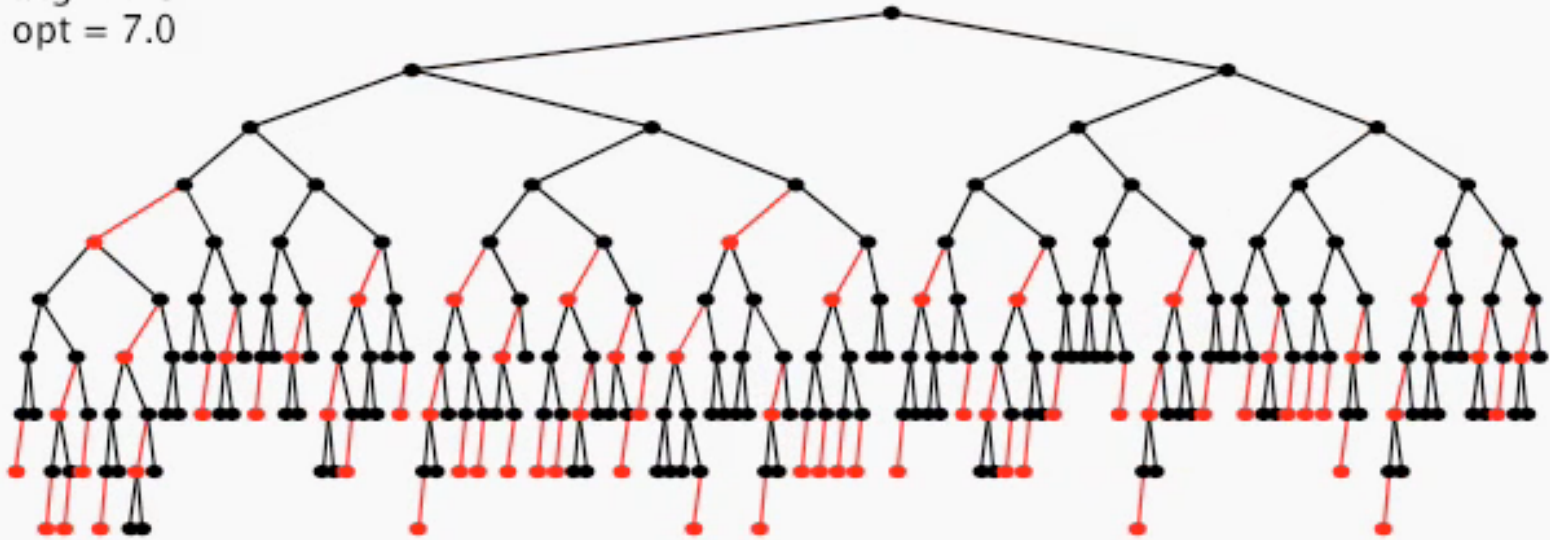
255 insertions in descending order

# Insertion in a LLRB tree: visualization



50 random insertions

## Insertion in a LLRB tree: visualization
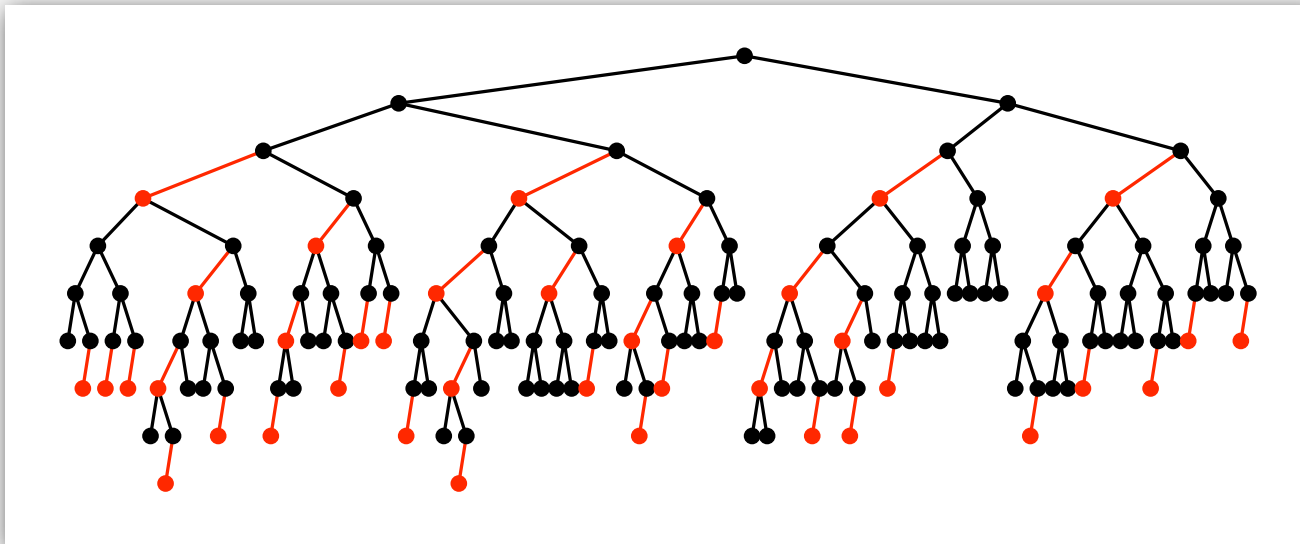


N = 255
max = 10
avg = 7.3
opt = 7.0

255 random insertions

## Balance in LLRB trees

Proposition.  Height of tree is ≤ 2 lg N in the worst case.

Pf.

- Every path from root to null link has same number of black links.
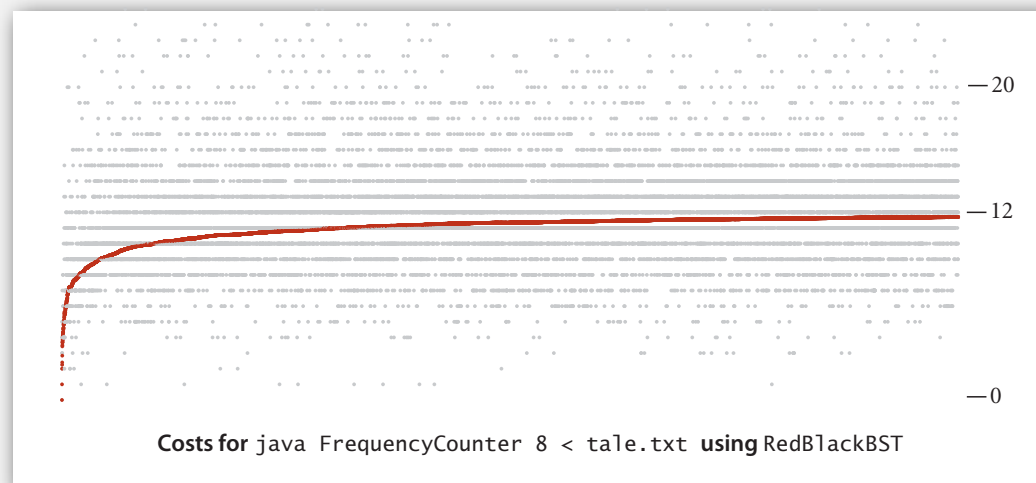- Never two red links in-a-row.



Property.  Height of tree is ~ 1.00 lg N in typical applications.

# ST implementations:  summary

| implementation | guarantee | | | average case | | | ordered iteration? | operations on keys |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | search | insert | delete | search hit | insert | delete | | |
| sequential search (linked list) | N | N | N | N/2 | N | N/2 | no | `equals()` |
| binary search (ordered array) | lg N | N | N | lg N | N/2 | N/2 | yes | `compareTo()` |
| BST | N | N | N | 1.39 lg N | 1.39 lg N | ? | yes | `compareTo()` |
| 2-3 tree | c lg N | c lg N | c lg N | c lg N | c lg N | c lg N | yes | `compareTo()` |
| red-black tree | 2 lg N | 2 lg N | 2 lg N | 1.00 lg N * | 1.00 lg N * | 1.00 lg N * | yes | `compareTo()` |

* exact value of coefficient unknown but extremely close to 1



**Costs for** `java FrequencyCounter 8 < tale.txt` **using** RedBlackBST

# Why left-leaning trees?

```
private Node put(Node x, Key key, Value val, boolean sw)
{
   if (x == null)
      return new Node(key, value, RED);
   int cmp = key.compareTo(x.key);

   if (isRed(x.left) && isRed(x.right))
   {
      x.color = RED;
      x.left.color  = BLACK;
      x.right.color = BLACK;
   }
   if (cmp < 0)
   {
      x.left = put(x.left, key, val, false);
      if (isRed(x) && isRed(x.left) && sw)
         x = rotateRight(x);
      if (isRed(x.left) && isRed(x.left.left))
      {
         x = rotateRight(x);
         x.color = BLACK; x.right.color = RED;
      }
   }
   else if (cmp > 0)
   {
      x.right = put(x.right, key, val, true);
      if (isRed(h) && isRed(x.right) && !sw)
         x = rotateLeft(x);
      if (isRed(h.right) && isRed(h.right.right))
      {
         x = rotateLeft(x);
         x.color = BLACK; x.left.color = RED;
      }
   }
   else x.val = val;
   return x;
}
```

```
public Node put(Node h, Key key, Value val)
{
   if (h == null)
      return new Node(key, val, RED);
   int cmp = kery.compareTo(h.key);
   if (cmp < 0)
      h.left  = put(h.left,  key, val);
   else if (cmp > 0)
      h.right = put(h.right, key, val);
   else h.val = val;

   if (isRed(h.right) && !isRed(h.left))
      h = rotateLeft(h);
   if (isRed(h.left) && isRed(h.left.left))
      h = rotateRight(h);
   if (isRed(h.left) && isRed(h.right))
      h = flipColors(h);

   return h;
}
```

straightforward
(if you've paid attention)

extremely tricky

## Why left-leaning trees?

**Simplified code.**
- Left-leaning restriction reduces number of cases.
- Short inner loop.

**Same ideas simplify implementation of other operations.**
- Delete min/max.
- Arbitrary delete.

*2008*
*1978*

**Improves widely-used algorithms.**
- AVL trees, 2-3 trees, 2-3-4 trees.
- Red-black trees.

*1972*

**Bottom line.** Left-leaning red-black trees are the simplest balanced BST to implement and the fastest in practice.

‣ 2-3-4 trees

‣ red-black trees

‣ **B-trees**

## File system model

Page.  Contiguous block of data (e.g., a file or 4096-byte chunk).

Probe.  First access to a page (e.g., from disk to memory).



slow                          fast

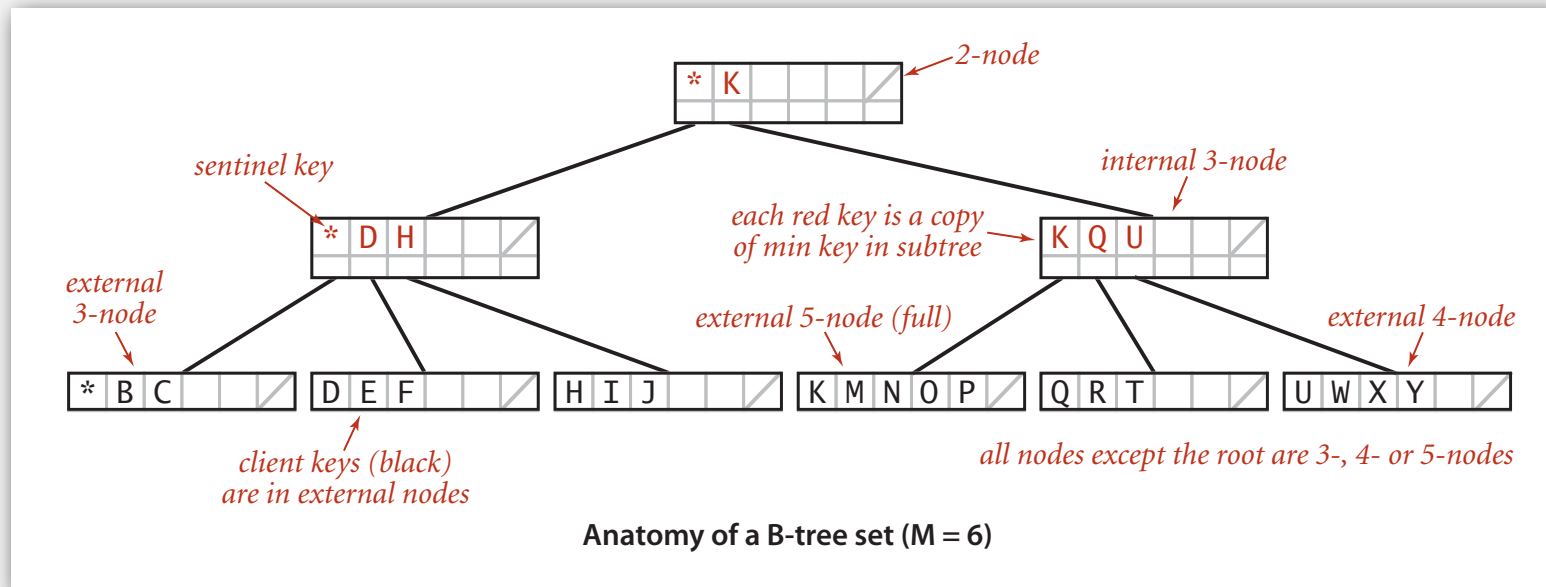Model.  Time required for a probe is much larger than time to access data within a page.

Goal.  Access data using minimum number of probes.

# B-trees (Bayer-McCreight, 1972)

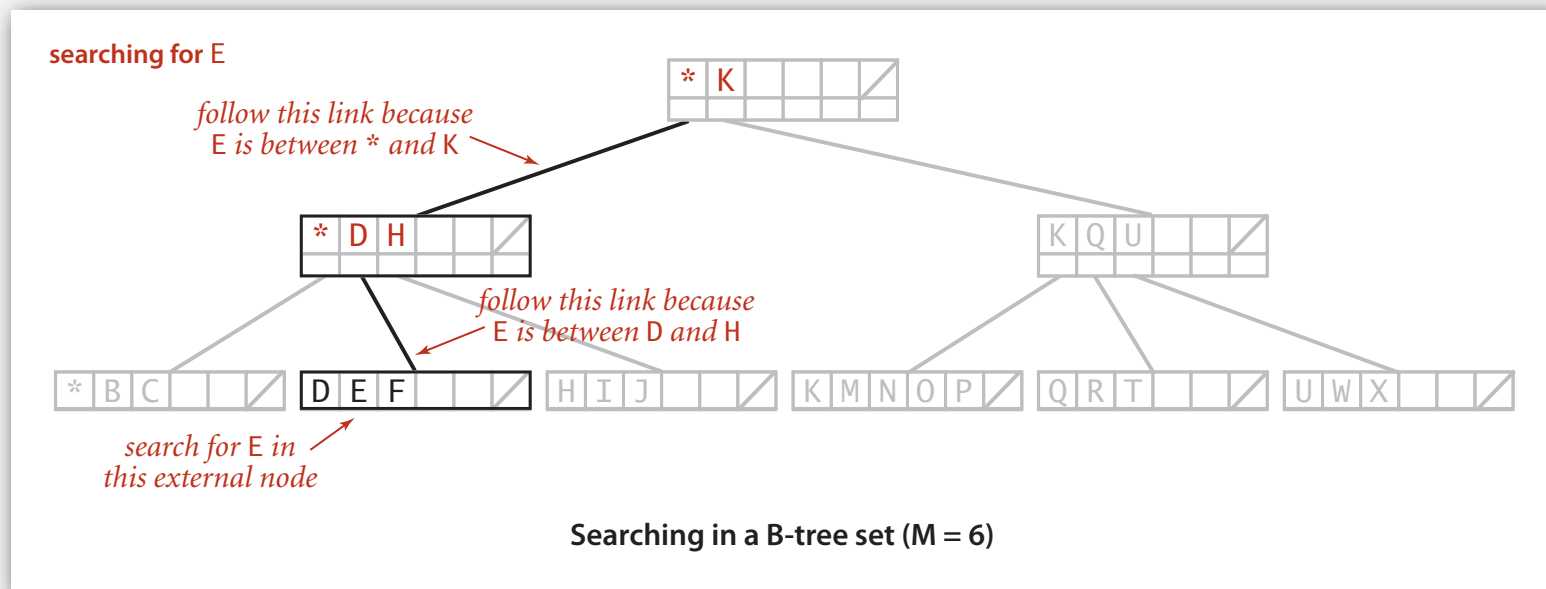**B-tree.** Generalize 2-3 trees by allowing up to M-1 key-link pairs per node.

- At least 2 key-link pairs at root.
- At least M/2 key-link pairs in other nodes.
- External nodes contain client keys.
- Internal nodes contain copies of keys to guide search.

*choose M as large as possible so that M links fit in a page, e.g., M = 1000*



*2-node*

*sentinel key*

*internal 3-node*

*each red key is a copy of min key in subtree*

*external 3-node*

*external 5-node (full)*

*external 4-node*

`* B C` `D E F` `H I J` `K M N O P` `Q R T` `U W X Y`

*client keys (black) are in external nodes*

*all nodes except the root are 3-, 4- or 5-nodes*
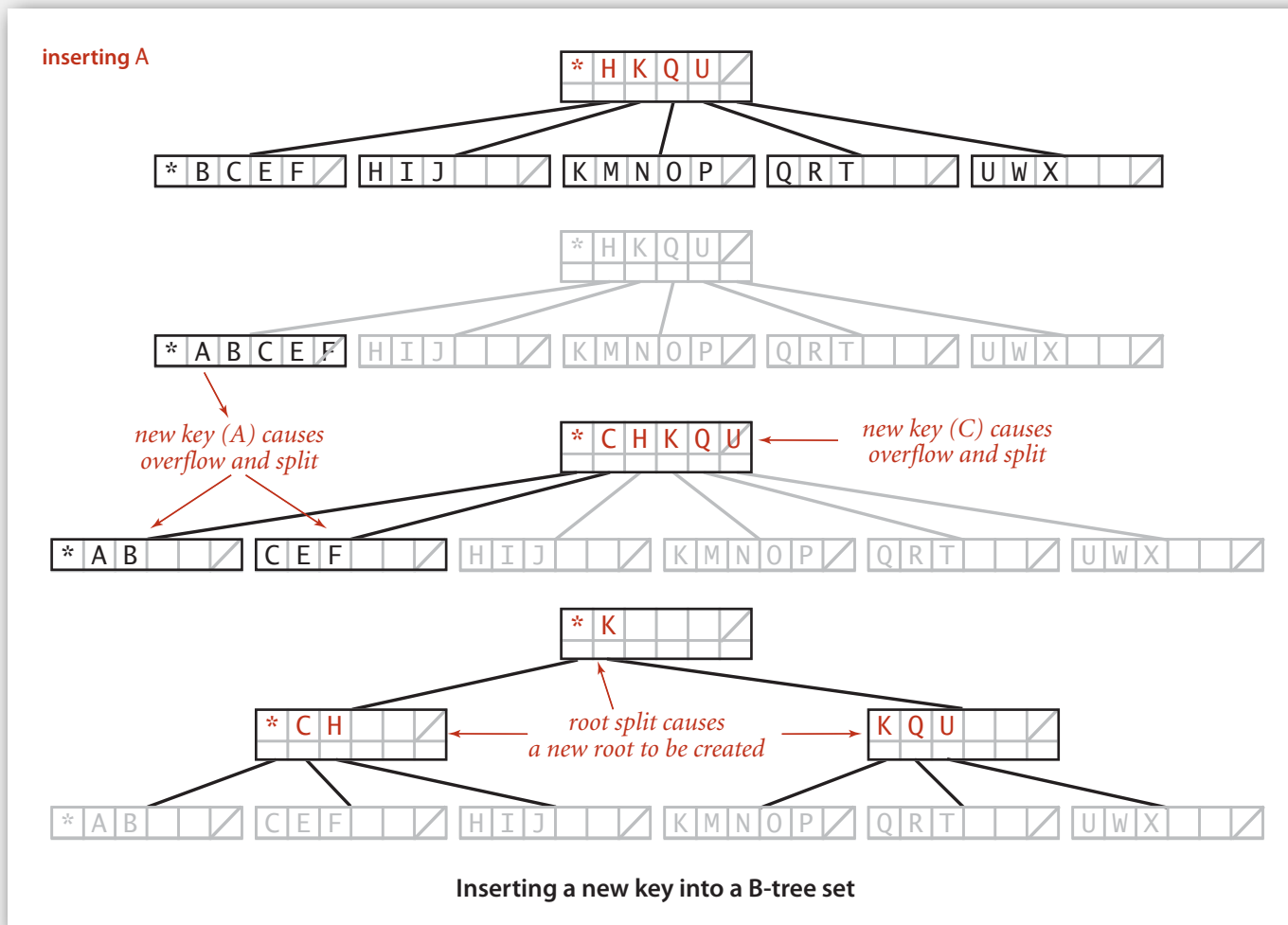
**Anatomy of a B-tree set (M = 6)**

# Searching in a B-tree

- Start at root.
- Find interval for search key and take corresponding link.
- Search terminates in external node.



**searching for** E

*follow this link because E is between * and K*

*follow this link because E is between D and H*

*search for E in this external node*

**Searching in a B-tree set (M = 6)**

# Insertion in a B-tree

- Search for new key.
- Insert at bottom.
- Split nodes with M key-link pairs on the way up the tree.



**inserting** A

*new key (A) causes overflow and split*

*new key (C) causes overflow and split*

*root split causes a new root to be created*

**Inserting a new key into a B-tree set**
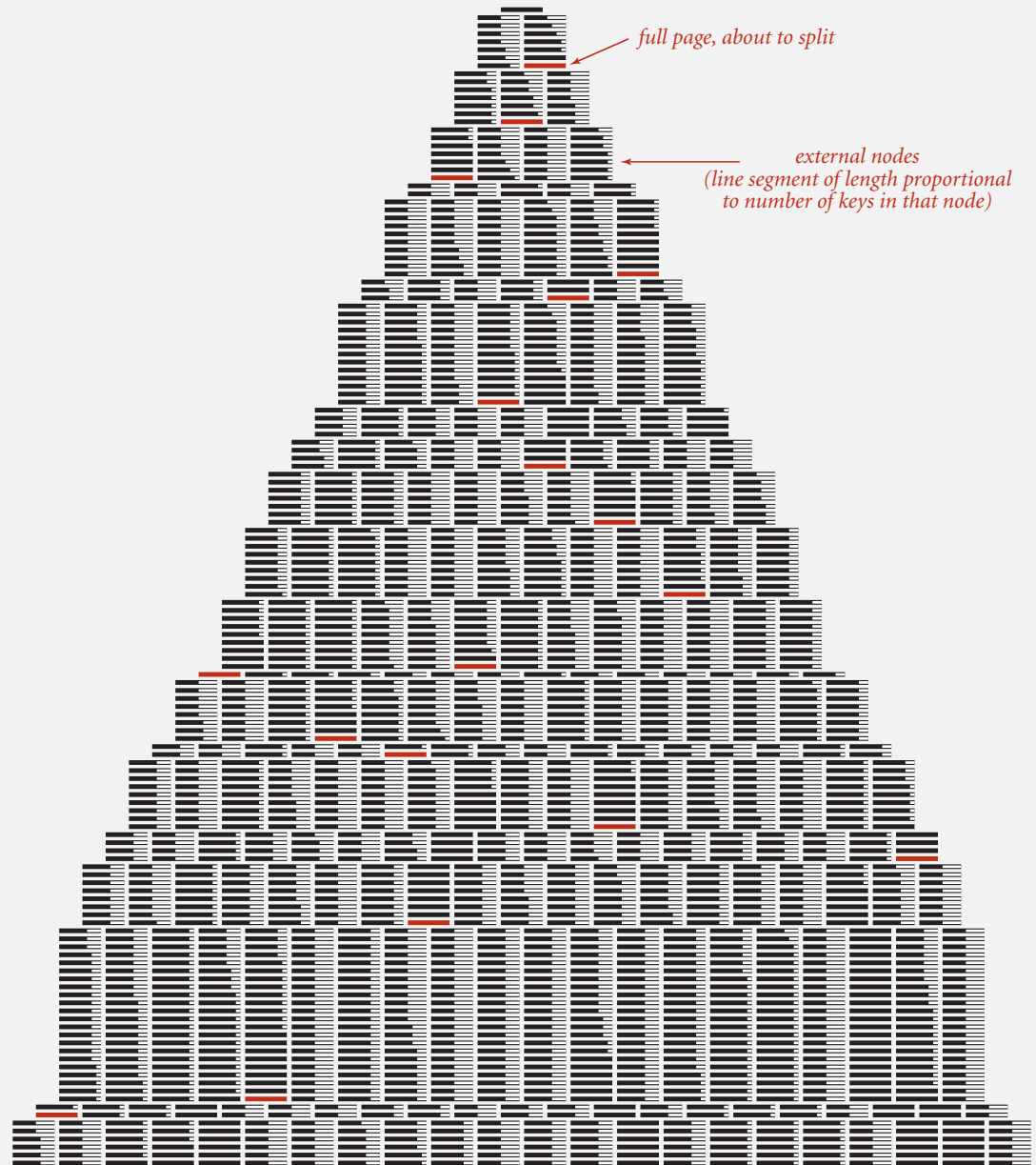
Balance in B-tree

Proposition.  A search or an insertion in a B-tree of order M with N keys requires between $\log_{M-1}N$ and $\log_{M/2}N$ probes.

Pf.  All internal nodes (besides root) have between M/2 and M-1 links.

In practice.  Number of probes is at most 4. $\longleftarrow$ $\quad$ M = 1000; N = 62 billion
$$\log_{M/2}N \leq 4$$

Optimization.  Always keep root page in memory.

# Building a large B tree



full page, about to split

external nodes
(line segment of length proportional
to number of keys in that node)

Balanced trees in the wild
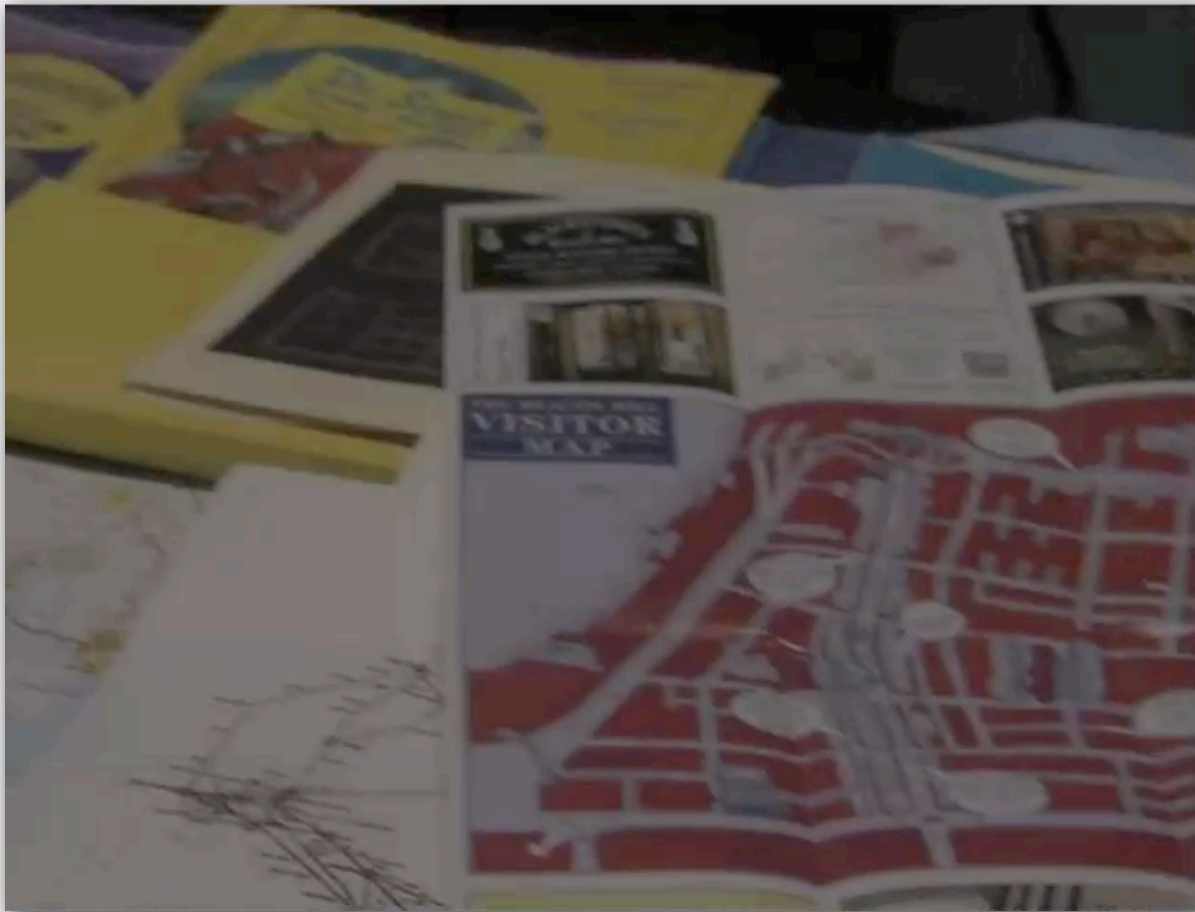
Red-black trees are widely used as system symbol tables.
- Java: `java.util.TreeMap`, `java.util.TreeSet`.
- C++ STL: map, multimap, multiset.
- Linux kernel: completely fair scheduler, `linux/rbtree.h`.

B-tree variants. B+ tree, B*tree, B# tree, …

B-trees (and variants) are widely used for file systems and databases.
- Windows: HPFS.
- Mac: HFS, HFS+.
- Linux: ReiserFS, XFS, Ext3FS, JFS.
- Databases: ORACLE, DB2, INGRES, SQL, PostgreSQL.

# Red-black trees in the wild





*Common sense. Sixth sense. Together they're the FBI's newest team.*

# Red-black trees in the wild

**ACT FOUR**

FADE IN:

48        INT. FBI HQ – NIGHT                                          48

Antonio is at THE COMPUTER as Jess explains herself to Nicole
and Pollock. The CONFERENCE TABLE is covered with OPEN
REFERENCE BOOKS, TOURIST GUIDES, MAPS and REAMS OF PRINTOUTS.

                     JESS
It was the red door again.

                     POLLOCK
I thought the red door was the storage
container.

                     JESS
But it wasn't red anymore.  It was
black.

                     ANTONIO
So red turning to black means...
what?

                     POLLOCK
Budget deficits?  Red ink, black
ink?

                     NICOLE
Yes.  I'm sure that's what it is.
But maybe we should come up with a
couple other options, just in case.

Antonio refers to his COMPUTER SCREEN, which is filled with
mathematical equations.

                     ANTONIO
It could be an algorithm from a binary
search tree.  A red-black tree tracks
every simple path from a node to a
descendant leaf with the same number
of black nodes.

                     JESS
Does that help you with girls?

Nicole is tapping away at a computer keyboard.  She finds
something.

53

# 3.4  Hash Tables



‣ **hash functions**

‣ **separate chaining**

‣ **linear probing**

‣ **applications**

# Optimize judiciously

> *"More computing sins are committed in the name of efficiency (without necessarily achieving it) than for any other single reason—including blind stupidity."* — William A. Wulf

> *"We should forget about small efficiencies, say about 97% of the time: premature optimization is the root of all evil."* — Donald E. Knuth

> *"We follow two rules in the matter of optimization:*
> *Rule 1: Don't do it.*
> *Rule 2 (for experts only). Don't do it yet - that is, not until you have a perfectly clear and unoptimized solution."* — M. A. Jackson

Reference: Effective Java by Joshua Bloch

## ST implementations: summary

| implementation | guarantee | | | average case | | | ordered iteration? | operations on keys |
|---|---|---|---|---|---|---|---|---|
| | search | insert | delete | search hit | insert | delete | | |
| sequential search (linked list) | N | N | N | N/2 | N | N/2 | no | `equals()` |
| binary search (ordered array) | lg N | N | N | lg N | N/2 | N/2 | yes | `compareTo()` |
| BST | N | N | N | 1.38 lg N | 1.38 lg N | ? | yes | `compareTo()` |
| red-black tree | 2 lg N | 2 lg N | 2 lg N | 1.00 lg N | 1.00 lg N | 1.00 lg N | yes | `compareTo()` |

Q.  Can we do better?

A.  Yes, but with different access to the data.

## Hashing:  basic plan

Save items in a key-indexed table (index is a function of the key).

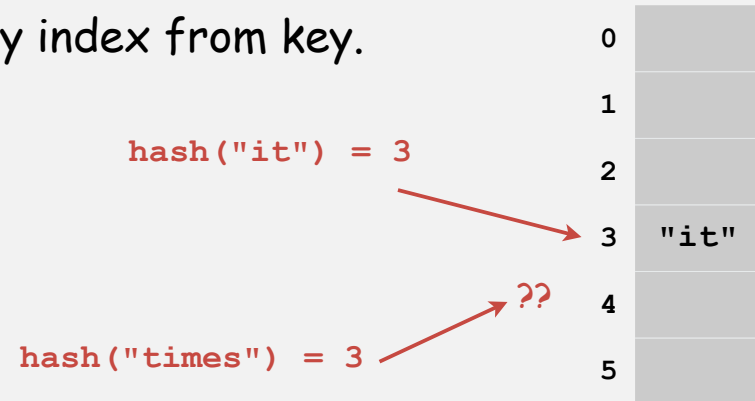Hash function.  Method for computing array index from key.

hash("it") = 3

Issues.

- Computing the hash function.
- Equality test:  Method for checking whether two keys are equal.

```
0
1
2
3   "it"
4
5
```

## Hashing:  basic plan

Save items in a key-indexed table (index is a function of the key).

Hash function.  Method for computing array index from key.

```
hash("it") = 3
```

```
hash("times") = 3
```

```
0
1
2
3   "it"
?? 4
5
```

Issues.

- Computing the hash function.
- Equality test:  Method for checking whether two keys are equal.
- Collision resolution:  Algorithm and data structure
  to handle two keys that hash to the same array index.

Classic space-time tradeoff.

- No space limitation:  trivial hash function with key as index.
- No time limitation:  trivial collision resolution with sequential search.
- Limitations on both time and space:  hashing (the real world).

## Equality test

Needed because hash methods do not use `compareTo()`.

All Java classes inherit a method `equals()`.

Java requirements. For any references `x`, `y` and `z`:

- Reflexive: `x.equals(x)` is `true`.
- Symmetric: `x.equals(y)` iff `y.equals(x)`.
- Transitive: if `x.equals(y)` and `y.equals(z)`, then `x.equals(z)`.
- Non-null: `x.equals(null)` is `false`.

equivalence
relation

do `x` and `y` refer to
the same object?

Default implementation. `(x == y)`

Customized implementations. `Integer`, `Double`, `String`, `File`, `URL`, `Date`, ...

User-defined implementations. Some care needed.

# Implementing equals for user-defined types

Seems easy

```java
public        class Record
{
    private final String name;
    private final long val;
    ...

    public boolean equals(Record y)
    {




        Record that =           y;
        return (this.val == that.val) &&
               (this.name.equals(that.name));
    }
}
```

check that all significant
fields are the same

# Implementing equals for user-defined types

Seems easy, but requires some care.

no safe way to use `equals()` with inheritance

```java
public final class Record
{
    private final String name;
    private final long val;
    ...

    public boolean equals(Object y)
    {
        if (y == this) return true;

        if (y == null) return false;

        if (y.getClass() != this.getClass())
            return false;

        Record that = (Record) y;
        return (this.val == that.val) &&
               (this.name.equals(that.name));
    }
}
```

must be `Object`.
Why? Experts still debate.

optimize for true object equality

check for `null`

objects must be in the same class

check that all significant fields are the same

## Computing the hash function

**Idealistic goal.** Scramble the keys uniformly to produce a table index.

- Efficiently computable.
- Each table index equally likely for each key.

*thoroughly researched problem,*
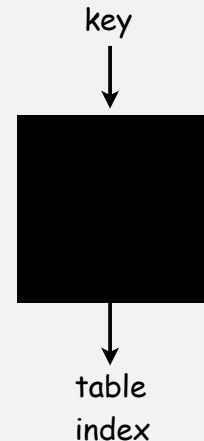*still problematic in practical applications*

key

↓

table
index

**Ex 1.** Phone numbers.

- Bad: first three digits.
- Better: last three digits.

**Ex 2.** Social Security numbers. ← *573 = California, 574 = Alaska*
*(assigned in chronological order within geographic region)*

- Bad: first three digits.
- Better: last three digits.

**Practical challenge.** Need different approach for each key type.
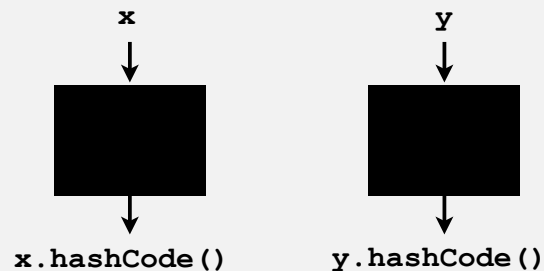
## Java's hash code conventions

All Java classes inherit a method `hashCode()`, which returns a 32-bit `int`.

Requirement. If `x.equals(y)`, then `(x.hashCode() == y.hashCode())`.

Highly desirable. If `!x.equals(y)`, then `(x.hashCode() != y.hashCode())`.



```
x                    y
↓                    ↓
■                    ■

↓                    ↓
x.hashCode()      y.hashCode()
```

Default implementation.  Memory address of `x`.
Customized implementations.  `Integer, Double, String, File, URL, Date, …`
User-defined types.  Users are on their own.

# Implementing hash code:  integers and doubles

```
public final class Integer
{
    private final int value;
    ...

    public int hashCode()
    {   return value;   }
}
```

```
public final class Double
{
    private final double value;
    ...

    public int hashCode()
    {
        long bits = doubleToLongBits(value);
        return (int) (bits ^ (bits >>> 32));
    }
}
```

convert to IEEE 64-bit representation;
xor most significant 32-bits
with least significant 32-bits

## Implementing hash code:  strings

```
public final class String
{
    private final char[] s;
    ...

    public int hashCode()
    {
        int hash = 0;
        for (int i = 0; i < length(); i++)
            hash = s[i] + (31 * hash);
        return hash;
    }
}
```

i<sup>th</sup> character of s

| char | Unicode |
|------|---------|
| … | … |
| 'a' | 97 |
| 'b' | 98 |
| 'c' | 99 |
| … | … |

- Horner's method to hash string of length L:  L multiplies/adds.
- Equivalent to  $h = 31^{L-1} \cdot s^0 + \ldots + 31^2 \cdot s^{L-3} + 31^1 \cdot s^{L-2} + 31^0 \cdot s^{L-1}$.

Ex.
```
String s = "call";
int code = s.hashCode();
```
$3045982 = 99 \cdot 31^3 + 97 \cdot 31^2 + 108 \cdot 31^1 + 108 \cdot 31^0$

$= 108 + 31 \cdot (108 + 31 \cdot (97 + 31 \cdot (99)))$

## A poor hash code

Ex. Strings (in Java 1.1).

- For long strings: only examine 8-9 evenly spaced characters.
- Benefit: saves time in performing arithmetic.

```
public int hashCode()
{
    int hash = 0;
    int skip = Math.max(1, length() / 8);
    for (int i = 0; i < length(); i += skip)
        hash = s[i] + (37 * hash);
    return hash;
}
```

- Downside: great potential for bad collision patterns.

```
http://www.cs.princeton.edu/introcs/13loop/Hello.java
http://www.cs.princeton.edu/introcs/13loop/Hello.class
http://www.cs.princeton.edu/introcs/13loop/Hello.html
http://www.cs.princeton.edu/introcs/13loop/index.html
http://www.cs.princeton.edu/introcs/12type/index.html
```

## Implementing hash code: user-defined types

```
public final class Record
{
   private String name;
   private int id;
   private double value;

   public Record(String name, int id, double value)
   {  /* as before */  }


   ...

   public boolean equals(Object y)
   {  /* as before */  }

   public int hashCode()
   {
      int hash = 17;
      hash = 31*hash + name.hashCode();
      hash = 31*hash + id;
      hash = 31*hash + Double.valueOf(value).hashCode();
      return hash;
   }
}
```

nonzero constant

typically a small prime

## Hash code design

"Standard" recipe for user-defined types.

- Combine each significant field using the 31x + y rule.
- If field is a primitive type, use built-in hash code.
- If field is an array, apply to each element.
- If field is an object, apply rule recursively.

In practice.   Recipe works reasonably well; used in Java libraries.

In theory.  Need a theorem for each type to ensure reliability.

Basic rule.  Need to use the whole key to compute hash code;

consult an expert for state-of-the-art hash codes.

## Modular hashing

Hash code.  An `int` between $-2^{31}$ and $2^{31}-1$.

Hash function.  An `int` between `0` and `M-1` (for use as array index).

typically a prime or power of 2

```
private int hash(Key key)
{   return key.hashCode() % M;   }
```

bug

```
private int hash(Key key)
{   return Math.abs(key.hashCode()) % M;   }
```

1-in-a-billion bug

```
private int hash(Key key)
{   return (key.hashCode() & 0x7fffffff) % M;   }
```

correct

## Uniform hashing assumption

Assumption J (uniform hashing hashing assumption).

Each key is equally likely to hash to an integer between 0 and M-1.

Bins and balls.  Throw balls uniformly at random into M bins.



```
0  1  2  3  4  5  6  7  8  9  10 11 12 13 14 15
```

Birthday problem.  Expect two balls in the same bin after $\sim \sqrt{\pi M / 2}$ tosses.

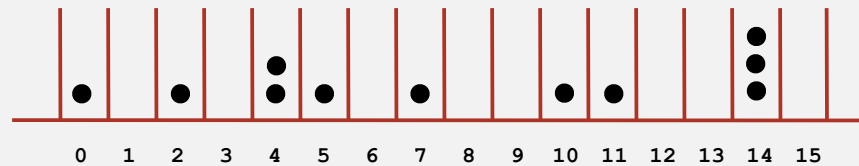Coupon collector.  Expect every bin has ≥ 1 ball after ~ M ln M tosses.

Load balancing.  After M tosses, expect most loaded bin has $\Theta(\log M / \log \log M)$ balls.

# Uniform hashing assumption

**Assumption J (uniform hashing hashing assumption).**

Each key is equally likely to hash to an integer between 0 and M-1.

**Bins and balls.** Throw balls uniformly at random into M bins.



Hash value frequencies for words in Tale of Two Cities (M = 97)

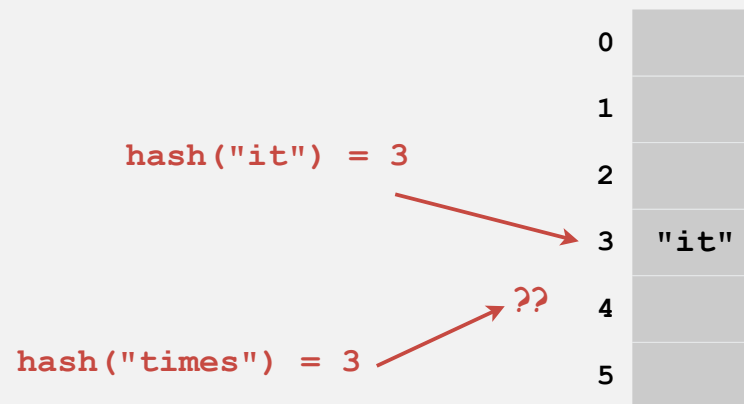Java's `string` data uniformly distribute the keys of Tale of Two Cities

▸ hash functions

▸ **separate chaining**

▸ linear probing

▸ applications

## Collisions

Collision.  Two distinct keys hashing to same index.

- Birthday problem $\Rightarrow$ can't avoid collisions unless you have a ridiculous amount (quadratic) of memory.
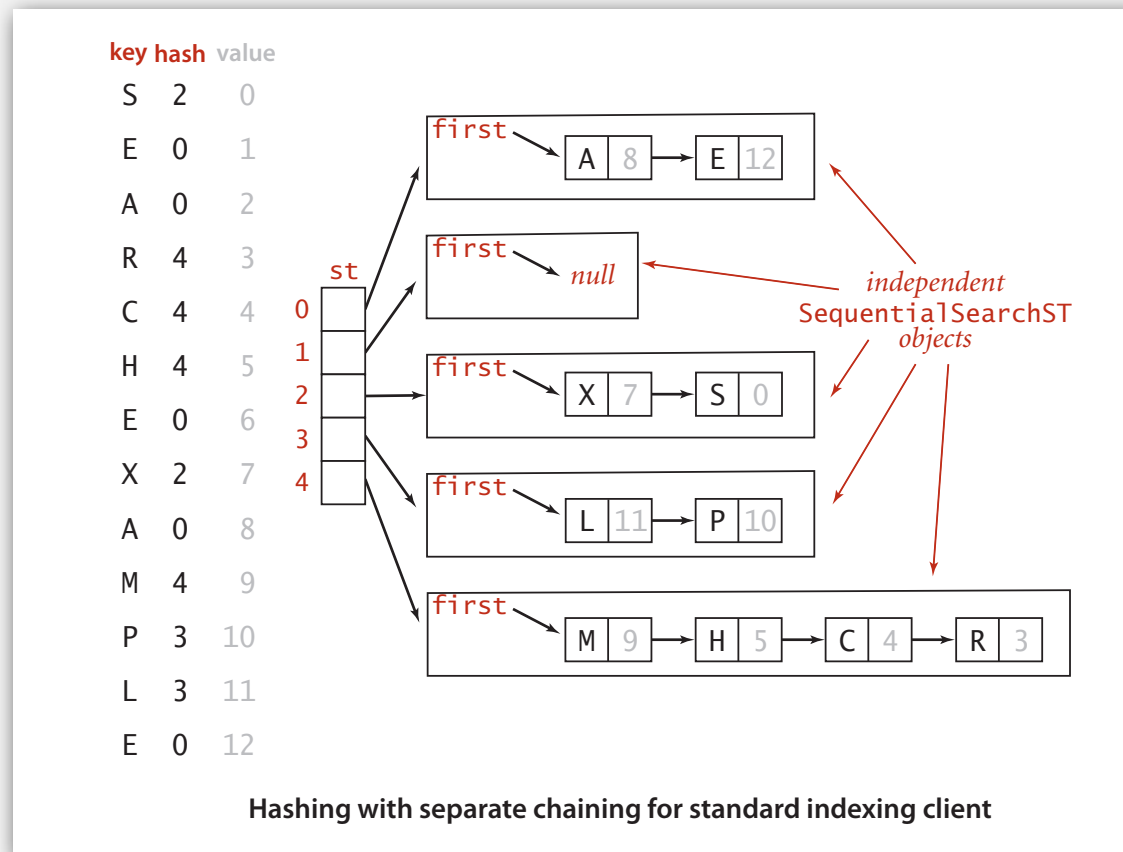- Coupon collector + load balancing $\Rightarrow$ collisions will be evenly distributed.

Challenge.  Deal with collisions efficiently.

hash("it") = 3

hash("times") = 3

??

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | |
| 3 | "it" |
| 4 | |
| 5 | |

# Separate chaining ST

Use an array of M < N linked lists.  [H. P. Luhn, IBM 1953]

- Hash:  map key to integer i between 0 and M-1.
- Insert:  put at front of $i^{th}$ chain (if not already there).
- Search:  only need to search $i^{th}$ chain.



**Hashing with separate chaining for standard indexing client**

## Separate chaining ST:  Java implementation

```
public class SeparateChainingHashST<Key, Value>
{
   private int N;       // number of key-value pairs
   private int M;       // hash table size
   private SequentialSearchST<Key, Value> [] st; // array of STs

   public SeparateChainingHashST()           ←——— array doubling code omitted
   {  this(997);   }

   public SeparateChainingHashST(int M)
   {
      this.M = M;
      st = (SequentialSearchST<Key, Value>[]) new SequentialSearchST[M];
      for (int i = 0; i < M; i++)
         st[i] = new SequentialSearchST<Key, Value>();
   }
   private int hash(Key key)
   {  return (key.hashCode() & 0x7fffffff) % M;   }

   public Value get(Key key)
   {  return st[hash(key)].get(key);   }

   public void put(Key key, Value val)
   {  st[hash(key)].put(key, val);   }
}
```
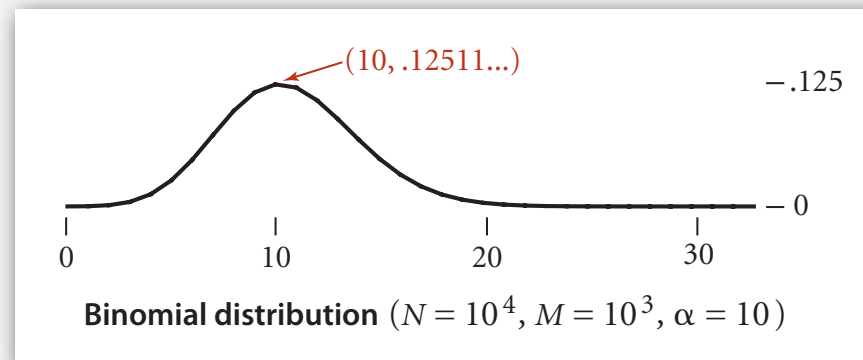
## Analysis of separate chaining

**Proposition K.** Under uniform hashing assumption, probability that the number of keys in a list is within a constant factor of N/M is extremely close to 1.

**Pf sketch.** Distribution of list size obeys a binomial distribution.

$(10, .12511...)$

$-.125$

$-0$

0   10   20   30

**Binomial distribution** $(N = 10^4, M = 10^3, \alpha = 10)$

`equals()` and `hashCode()`

**Consequence.** Number of probes for search/insert is proportional to N/M.
- M too large $\Rightarrow$ too many empty chains.
- M too small $\Rightarrow$ chains too long.
- Typical choice: M ~ N/5 $\Rightarrow$ constant-time ops.

M times faster than sequential search

# Collision resolution:  open addressing

Open addressing.  [Amdahl-Boehme-Rocherster-Samuel, IBM 1953]

When a new key collides, find next empty slot, and put it there.

| | |
|---|---|
| st[0] | jocularly |
| st[1] | *null* |
| st[2] | listen |
| st[3] | suburban |
| ⋮ | *null* |
| st[30000] | browsing |

linear probing (M = 30001, N = 15000)

## Linear probing

Use an array of size M > N.

- Hash:  map key to integer i between 0 and M-1.
- Insert:  put at table index i if free; if not try i+1, i+2, etc.
- Search:  search table index i; if occupied but no match, try i+1, i+2, etc.

| - | - | - | S | H | - | - | A | C | E | R | - | - |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |

| - | - | - | S | H | - | - | A | C | E | R | I | - |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |

insert I
hash(I) = 11

| - | - | - | S | H | - | - | A | C | E | R | I | N |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |

insert N
hash(N) = 8

# Linear probing:  trace of standard indexing client

| key | hash | value | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| S | 6 | 0 | | | | | | | | S 0 | | | | | | | | | |
| E | 10 | 1 | | | | | | | | S 0 | | | | E 1 | | | | | |
| A | 4 | 2 | | | | | | A 2 | | S 0 | | | | E 1 | | | | | |
| R | 14 | 3 | | | | | | A 2 | | S 0 | | | | E 1 | | | | R 3 | |
| C | 5 | 4 | | | | | | A 2 | C 5 | S 0 | | | | E 1 | | | | R 3 | |
| H | 4 | 5 | | | | | | A 2 | C 5 | S 0 | H 5 | | | E 1 | | | | R 3 | |
| E | 10 | 6 | | | | | | A 2 | C 5 | S 0 | H 5 | | | E (6) | | | | R 3 | |
| X | 15 | 7 | | | | | | A 2 | C 5 | S 0 | H 5 | | | E 6 | | | | R 3 | X 7 |
| A | 4 | 8 | | | | | | A (8) | C 5 | S 0 | H 5 | | | E 6 | | | | R 3 | X 7 |
| M | 1 | 9 | | | M 9 | | | A 8 | C 5 | S 0 | H 5 | | | E 6 | | | | R 3 | X 7 |
| P | 14 | 10 | | P 10 | M 9 | | | A 8 | C 5 | S 0 | H 5 | | | E 6 | | | | R 3 | X 7 |
| L | 6 | 11 | | P 10 | M 9 | | | A 8 | C 5 | S 0 | H 5 | L 11 | | E 6 | | | | R 3 | X 7 |
| E | 10 | 12 | | P 10 | M 9 | | | A 8 | C 5 | S 0 | H 5 | L 11 | | E (12) | | | | R 3 | X 7 |

*entries in red are new*

*entries in gray are untouched*

*keys in black are probes*

*probe sequence wraps to 0*

← keys[]
← vals[]

28

## Linear probing ST implementation

```
public class LinearProbingHashST<Key, Value>
{
    private int M = 30001;
    private Value[] vals = (Value[]) new Object[M];
    private Key[]   keys = (Key[])   new Object[M];

    private int hash(Key key) {  /* as before */  }

    public void put(Key key, Value val)
    {
        int i;
        for (i = hash(key); keys[i] != null; i = (i+1) % M)
            if (keys[i].equals(key))
                break;
        keys[i] = key;
        vals[i] = val;
    }

    public Value get(Key key)
    {
        for (int i = hash(key); keys[i] != null; i = (i+1) % M)
            if (key.equals(keys[i]))
                return vals[i];
        return null;
    }
}
```
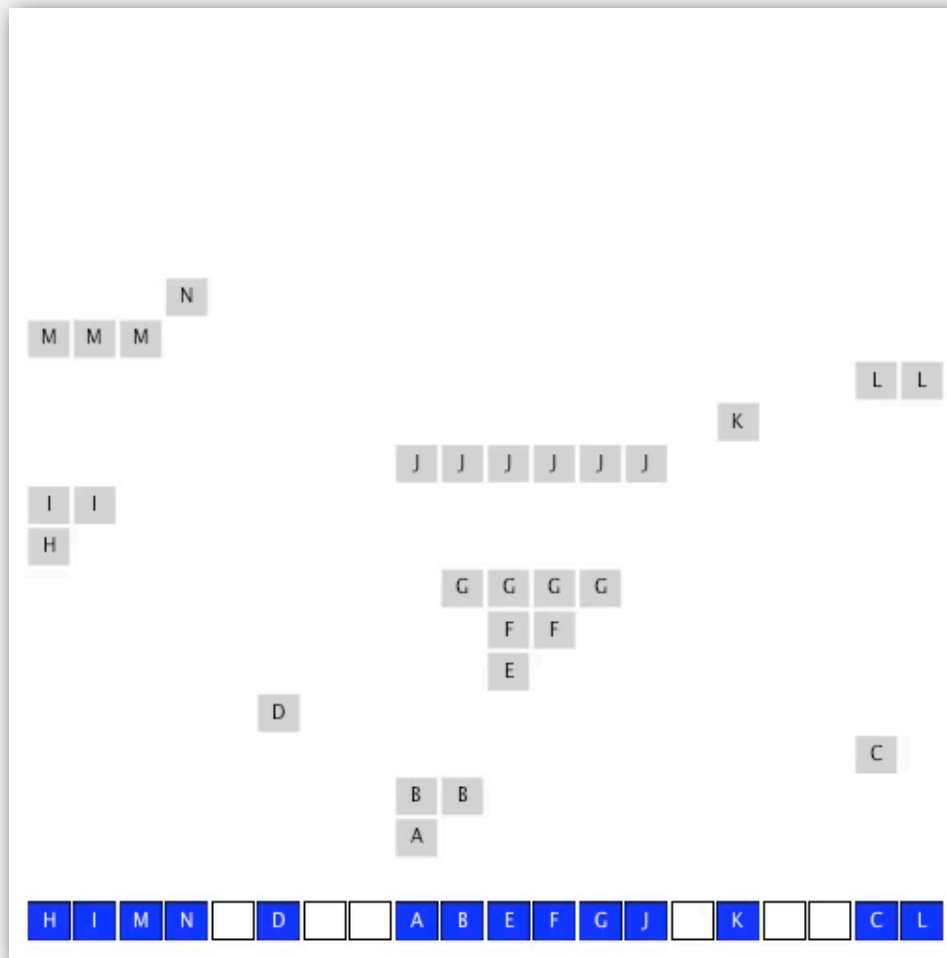
array doubling
code omitted

# Clustering

Cluster.  A contiguous block of items.

Observation.  New keys likely to hash into middle of big clusters.

# Knuth's parking problem

**Model.** Cars arrive at one-way street with M parking spaces.
Each desires a random space i:  if space i is taken, try i+1, i+2, …

**Q.** What is mean displacement of a car?



displacement = 3

**Empty.** With M/2 cars, mean displacement is ~ 3/2.
**Full.**   With M cars, mean displacement is ~ $\sqrt{\pi M / 8}$

# Analysis of linear probing

Proposition M. Under uniform hashing assumption, the average number of probes in a hash table of size M that contains N = $\alpha$ M keys is:

$$\sim \frac{1}{2}\left(1 + \frac{1}{1-\alpha}\right) \qquad \sim \frac{1}{2}\left(1 + \frac{1}{(1-\alpha)^2}\right)$$

<p style="text-align:center; color:red">search hit          search miss / insert</p>

Pf. [Knuth 1962]  A landmark in analysis of algorithms.

Parameters.

- M too large $\Rightarrow$ too many empty array entries.
- M too small $\Rightarrow$ search time blows up.
- Typical choice: $\alpha$ = N/M $\sim \frac{1}{2}$.

# probes for search hit is about 3/2
# probes for search miss is about 5/2

# ST implementations:  summary

| implementation | guarantee | | | average case | | | ordered iteration? | operations on keys |
|---|---|---|---|---|---|---|---|---|
| | search | insert | delete | search hit | insert | delete | | |
| sequential search (linked list) | N | N | N | N/2 | N | N/2 | no | `equals()` |
| binary search (ordered array) | lg N | N | N | lg N | N/2 | N/2 | yes | `compareTo()` |
| BST | N | N | N | 1.38 lg N | 1.38 lg N | ? | yes | `compareTo()` |
| red-black tree | 2 lg N | 2 lg N | 2 lg N | 1.00 lg N | 1.00 lg N | 1.00 lg N | yes | `compareTo()` |
| hashing | lg N * | lg N * | lg N * | 3-5 * | 3-5 * | 3-5 * | no | `equals()` |

\* under uniform hashing assumption

## Algorithmic complexity attacks

Q. Is the uniform hashing assumption important in practice?

A. Obvious situations: aircraft control, nuclear reactor, pacemaker.

A. Surprising situations: denial-of-service attacks.



malicious adversary learns your hash function
(e.g., by reading Java API) and causes a big pile-up
in single slot that grinds performance to a halt

Real-world exploits. [Crosby-Wallach 2003]

• Bro server: send carefully chosen packets to DOS the server,
  using less bandwidth than a dial-up modem.
• Perl 5.8.0: insert carefully chosen strings into associative array.
• Linux 2.4.20 kernel: save files with carefully chosen names.

# Algorithmic complexity attack on Java

Goal. Find family of strings with the same hash code.

Solution. The base-31 hash code is part of Java's string API.

| key | hashCode() |
|------|-----------|
| "Aa" | 2112 |
| "BB" | 2112 |

| key | hashCode() |
|------------|-------------|
| "AaAaAaAa" | −540425984 |
| "AaAaAaBB" | −540425984 |
| "AaAaBBAa" | −540425984 |
| "AaAaBBBB" | −540425984 |
| "AaBBAaAa" | −540425984 |
| "AaBBAaBB" | −540425984 |
| "AaBBBBAa" | −540425984 |
| "AaBBBBBB" | −540425984 |

| key | hashCode() |
|------------|-------------|
| "BBAaAaAa" | −540425984 |
| "BBAaAaBB" | −540425984 |
| "BBAaBBAa" | −540425984 |
| "BBAaBBBB" | −540425984 |
| "BBBBAaAa" | −540425984 |
| "BBBBAaBB" | −540425984 |
| "BBBBBBAa" | −540425984 |
| "BBBBBBBB" | −540425984 |

$2^N$ strings of length 2N that hash to same value!

## Diversion: one-way hash functions

**One-way hash function.** Hard to find a key that will hash to a desired value, or to find two keys that hash to same value.

**Ex.** MD4, MD5, SHA-0, SHA-1, SHA-2, WHIRLPOOL, RIPEMD-160.

known to be insecure

```
String password = args[0];
MessageDigest sha1 = MessageDigest.getInstance("SHA1");
byte[] bytes = sha1.digest(password);

/* prints bytes as hex string */
```

**Applications.** Digital fingerprint, message digest, storing passwords.

**Caveat.** Too expensive for use in ST implementations.

# Separate chaining vs. linear probing

## Separate chaining.

- Easier to implement delete.
- Performance degrades gracefully.
- Clustering less sensitive to poorly-designed hash function.

## Linear probing.

- Less wasted space.
- Better cache performance.

## Hashing: variations on the theme

Many improved versions have been studied.

**Two-probe hashing.** (separate chaining variant)
- Hash to two positions, put key in shorter of the two chains.
- Reduces average length of the longest chain to log log N.

**Double hashing.** (linear probing variant)
- Use linear probing, but skip a variable amount, not just 1 each time.
- Effectively eliminates clustering.
- Can allow table to become nearly full.

## Hashing vs. balanced trees

**Hashing.**

- Simpler to code.
- No effective alternative for unordered keys.
- Faster for simple keys (a few arithmetic ops versus log N compares).
- Better system support in Java for strings (e.g., cached hash code).

**Balanced trees.**

- Stronger performance guarantee.
- Support for ordered ST operations.
- Easier to implement `compareTo()` correctly than `equals()` and `hashCode()`.

**Java system includes both.**

- Red-black trees: `java.util.TreeMap, java.util.TreeSet.`
- Hashing: `java.util.HashMap, java.util.IdentityHashMap.`

# 3.5 Symbol Tables Applications

▸ sets
▸ dictionary clients
▸ indexing clients
▸ sparse vectors

‣ **sets**

‣ dictionary clients

‣ indexing clients

‣ sparse vectors

**Mathematical set.** A collection of distinct keys.

| | | |
|---|---|---|
| `public class SET<Key extends Comparable<Key>>` | | |
| `SET()` | *create an empty set* | |
| `void` | `add(Key key)` | *add the key to the set* |
| `boolean` | `contains(Key key)` | *is the key in the set?* |
| `void` | `remove(Key key)` | *remove the key from the set* |
| `int` | `size()` | *return the number of keys in the set* |
| `Iterator<Key>` | `iterator()` | *iterator through keys in the set* |

**Q.** How to implement?

## Exception filter

- Read in a list of words from one file.
- Print out all words from standard input that are { in, not in } the list.

```
% more list.txt
was it the of

% java WhiteList list.txt < tinyTale.txt
it was the of it was the of
it was the of it was the of
it was the of it was the of
it was the of it was the of
it was the of it was the of

% java BlackList list.txt < tinyTale.txt
best times worst times
age wisdom age foolishness
epoch belief epoch incredulity
season light season darkness
spring hope winter despair
```

list of exceptional words

## Exception filter applications

- Read in a list of words from one file.
- Print out all words from standard input that are { in, not in } the list.

| application | purpose | key | in list |
|---|---|---|---|
| spell checker | identify misspelled words | word | dictionary words |
| browser | mark visited pages | URL | visited pages |
| parental controls | block sites | URL | bad sites |
| chess | detect draw | board | positions |
| spam filter | eliminate spam | IP address | spam addresses |
| credit cards | check for stolen cards | number | stolen cards |

## Exception filter:  Java implementation

- Read in a list of words from one file.
- Print out all words from standard input that are { in, not in } the list.

```
public class WhiteList
{
    public static void main(String[] args)
    {
        SET<String> set = new SET<String>();          create empty set of strings

        In in = new In(args[0]);
        while (!in.isEmpty())                          read in whitelist
            set.add(in.readString());

        while (!StdIn.isEmpty())
        {
            String word = StdIn.readString();
            if (set.contains(word))                    print words in list
                StdOut.println(word);
        }
    }
}
```

## Exception filter: Java implementation

- Read in a list of words from one file.
- Print out all words from standard input that are { in, not in } the list.

```java
public class BlackList
{
    public static void main(String[] args)
    {
        SET<String> set = new SET<String>();          ← create empty set of strings

        In in = new In(args[0]);
        while (!in.isEmpty())                          ← read in blacklist
            set.add(in.readString());

        while (!StdIn.isEmpty())
        {
            String word = StdIn.readString();
            if (!set.contains(word))                   ← print words not in list
                StdOut.println(word);
        }
    }
}
```

▸ sets

▸ **dictionary clients**

▸ indexing clients

▸ sparse vectors

# Dictionary lookup

## Command-line arguments.

- A comma-separated value (CSV) file.
- Key field.
- Value field.

Ex 1. DNS lookup.

URL is key    IP is value

```
%  java LookupCSV ip.csv 0 1
adobe.com
192.150.18.60
www.princeton.edu
128.112.128.15
ebay.edu
Not found
```

IP is key    URL is value

```
%  java LookupCSV ip.csv 1 0
128.112.128.15
www.princeton.edu
999.999.999.99
Not found
```

```
% more ip.csv
www.princeton.edu,128.112.128.15
www.cs.princeton.edu,128.112.136.35
www.math.princeton.edu,128.112.18.11
www.cs.harvard.edu,140.247.50.127
www.harvard.edu,128.103.60.24
www.yale.edu,130.132.51.8
www.econ.yale.edu,128.36.236.74
www.cs.yale.edu,128.36.229.30
espn.com,199.181.135.201
yahoo.com,66.94.234.13
msn.com,207.68.172.246
google.com,64.233.167.99
baidu.com,202.108.22.33
yahoo.co.jp,202.93.91.141
sina.com.cn,202.108.33.32
ebay.com,66.135.192.87
adobe.com,192.150.18.60
163.com,220.181.29.154
passport.net,65.54.179.226
tom.com,61.135.158.237
nate.com,203.226.253.11
cnn.com,64.236.16.20
daum.net,211.115.77.211
blogger.com,66.102.15.100
fastclick.com,205.180.86.4
wikipedia.org,66.230.200.100
rakuten.co.jp,202.72.51.22
...
```

## Dictionary lookup

**Command-line arguments.**

- A comma-separated value (CSV) file.
- Key field.
- Value field.

**Ex 2.** Amino acids.

```
% java Lookup amino.csv 0 3
ACT
Threonine
TAG
Stop
CAT
Histidine
```

codon is key    name is value

```
% more amino.csv
TTT,Phe,F,Phenylalanine
TTC,Phe,F,Phenylalanine
TTA,Leu,L,Leucine
TTG,Leu,L,Leucine
TCT,Ser,S,Serine
TCC,Ser,S,Serine
TCA,Ser,S,Serine
TCG,Ser,S,Serine
TAT,Tyr,Y,Tyrosine
TAC,Tyr,Y,Tyrosine
TAA,Stop,Stop,Stop
TAG,Stop,Stop,Stop
TGT,Cys,C,Cysteine
TGC,Cys,C,Cysteine
TGA,Stop,Stop,Stop
TGG,Trp,W,Tryptophan
CTT,Leu,L,Leucine
CTC,Leu,L,Leucine
CTA,Leu,L,Leucine
CTG,Leu,L,Leucine
CCT,Pro,P,Proline
CCC,Pro,P,Proline
CCA,Pro,P,Proline
CCG,Pro,P,Proline
CAT,His,H,Histidine
CAC,His,H,Histidine
CAA,Gln,Q,Glutamine
CAG,Gln,Q,Glutamine
CGT,Arg,R,Arginine
CGC,Arg,R,Arginine
...
```

## Dictionary lookup

### Command-line arguments.

- A comma-separated value (CSV) file.
- Key field.
- Value field.

### Ex 3. Class list.

login is key
first name is value

```
% java Lookup classlist.csv 4 1
eberl
Ethan
nwebb
Natalie
```

login is key
precept is value

```
% java Lookup classlist.csv 4 3
dpan
P01
```

```
% more classlist.csv
13,Berl,Ethan Michael,P01,eberl
11,Bourque,Alexander Joseph,P01,abourque
12,Cao,Phillips Minghua,P01,pcao
11,Chehoud,Christel,P01,cchehoud
10,Douglas,Malia Morioka,P01,malia
12,Haddock,Sara Lynn,P01,shaddock
12,Hantman,Nicole Samantha,P01,nhantman
11,Hesterberg,Adam Classen,P01,ahesterb
13,Hwang,Roland Lee,P01,rhwang
13,Hyde,Gregory Thomas,P01,ghyde
13,Kim,Hyunmoon,P01,hktwo
11,Kleinfeld,Ivan Maximillian,P01,ikleinfe
12,Korac,Damjan,P01,dkorac
11,MacDonald,Graham David,P01,gmacdona
10,Michal,Brian Thomas,P01,bmichal
12,Nam,Seung Hyeon,P01,seungnam
11,Nastasescu,Maria Monica,P01,mnastase
11,Pan,Di,P01,dpan
12,Partridge,Brenton Alan,P01,bpartrid
13,Rilee,Alexander,P01,arilee
13,Roopakalu,Ajay,P01,aroopaka
11,Sheng,Ben C,P01,bsheng
12,Webb,Natalie Sue,P01,nwebb
...
```

## Dictionary lookup:  Java implementation

```
public class LookupCSV
{
   public static void main(String[] args)
   {
      In in = new In(args[0]);
      int keyField = Integer.parseInt(args[1]);
      int valField = Integer.parseInt(args[2]);

      ST<String, String> st = new ST<String, String>();
      while (!in.isEmpty())
      {
         String line = in.readLine();
         String[] tokens = database[i].split(",");
         String key = tokens[keyField];
         String val = tokens[valField];
         st.put(key, val);
      }

      while (!StdIn.isEmpty())
      {
         String s = StdIn.readString();
         if (!st.contains(s)) StdOut.println("Not found");
         else                 StdOut.println(st.get(s));
      }
   }
}
```

process input file

build symbol table

process lookups
with standard I/O

# File indexing

**Goal.** Index a PC (or the web).

# File indexing

**Goal.** Given a list of files specified as command-line arguments, create an index so that can efficiently find all files containing a given query string.

```
% ls *.txt
aesop.txt magna.txt moby.txt
sawyer.txt tale.txt

% java FileIndex *.txt
freedom
magna.txt moby.txt tale.txt

whale
moby.txt

lamb
sawyer.txt aesop.txt
```

```
% ls *.java

% java FileIndex *.java
BlackList.java Concordance.java
DeDup.java FileIndex.java ST.java
SET.java WhiteList.java

import
FileIndex.java SET.java ST.java

Comparator
null
```

**Solution.** Key = query string; value = set of files containing that string.

# File indexing

```
public class FileIndex
{
   public static void main(String[] args)
   {
      ST<String, SET<File>> st = new ST<String, SET<File>>();          symbol table

      for (String filename : args) {                                   list of file names
         File file = new File(filename);                               from command line
         In in = new In(file);
         while !(in.isEmpty())
         {
            String word = in.readString();                            for each word in file,
            if (!st.contains(word))                                   add file to
               st.put(s, new SET<File>());                            corresponding set
            SET<File> set = st.get(key);
            set.add(file);
         }
      }

      while (!StdIn.isEmpty())
      {
         String query = StdIn.readString();                           process queries
         StdOut.println(st.get(query));
      }
   }
}
```

# Book index

Goal.  Index for an e-book.

## Index

Abstract data type (ADT), 127-195
  abstract classes, 163
  classes, 129-136
  collections of items, 137-139
  creating, 157-164
  defined, 128
  duplicate items, 173-176
  equivalence-relations, 159-162
  FIFO queues, 165-171
  first-class, 177-186
  generic operations, 273
  index items, 177
  *insert/remove* operations, 138-139
  modular programming, 135
  polynomial, 188-192
  priority queues, 375-376
  pushdown stack, 138-156
  stubs, 135
  symbol table, 497-506
ADT interfaces
  array (myArray), 274
  complex number (Complex), 181
  existence table (ET), 663
  full priority queue (PQfull), 397
  indirect priority queue (PQi), 403
  item (myItem), 273, 498
  key (myKey), 498
  polynomial (Poly), 189
  point (Point), 134
  priority queue (PQ), 375
  queue of int (intQueue), 166

  stack of int (intStack), 140
  symbol table (ST), 503
  text index (TI), 525
  union–find (UF), 159
Abstract in-place merging, 351-353
Abstract operation, 10
Access control state, 131
Actual data, 31
Adapter class, 155-157
Adaptive sort, 268
Address, 84-85
Adjacency list, 120-123
  depth-first search, 251-256
Adjacency matrix, 120-122
Ajtai, M., 464
Algorithm, 4-6, 27-64
  abstract operations, 10, 31, 34-35
  analysis of, 6
  average-/worst-case performance, 35, 60-62
  big-Oh notation, 44-47
  binary search, 56-59
  computational complexity, 62-64
  efficiency, 6, 30, 32
  empirical analysis, 30-32, 58
  exponential-time, 219
  implementation, 28-30
  logarithm function, 40-43
  mathematical analysis, 33-36, 58
  primary parameter, 36
  probabilistic, 331
  recurrences, 49-52, 57
  recursive, 198
  running time, 34-40
  search, 53-56, 498
  steps in, 22-23
  *See also* Randomized algorithm
Amortization approach, 557, 627
Arithmetic operator, 177-179, 188, 191
Array, 12, 83
  binary search, 57
  dynamic allocation, 87

  and linked lists, 92, 94-95
  merging, 349-350
  multidimensional, 117-118
  references, 86-87, 89
  sorting, 265-267, 273-276
  and strings, 119
  two-dimensional, 117-118, 120-124
  vectors, 87
  visualizations, 295
  *See also* Index, array
Array representation
  binary tree, 381
  FIFO queue, 168-169
  linked lists, 110
  polynomial ADT, 191-192
  priority queue, 377-378, 403, 406
  pushdown stack, 148-150
  random queue, 170
  symbol table, 508, 511-512, 521
Asymptotic expression, 45-46
Average deviation, 80-81
Average-case performance, 35, 60-61
AVL tree, 583

B tree, 584, 692-704
  external/internal pages, 695
  4-5-6-7-8 tree, 693-704
  Markov chain, 701
  *remove*, 701-703
  *search/insert*, 697-701
  *select/sort*, 701
Balanced tree, 238, 555-598
  B tree, 584
  bottom-up, 576, 584-585
  height-balanced, 583
  indexed sequential access, 690-692
  performance, 575-576, 581-582, 595-598
  randomized, 559-564
  red–black, 577-585
  skip lists, 587-594
  splay, 566-571

727

## Concordance

Goal. Preprocess a text corpus to support concordance queries: given a word, find all occurrences with their immediate contexts.

```
% java Concordance tale.txt
cities
tongues of the two *cities* that were blended in

majesty
their turnkeys and the *majesty* of the law fired
me treason against the *majesty* of the people in
of his most gracious *majesty* king george the third

princeton
no matches
```

## Concordance

```java
public class Concordance
{
   public static void main(String[] args)
   {
      In in = new In(args[0]);
      String[] words = StdIn.readAll().split("\\s+");
      ST<String, SET<Integer>> st = new ST<String, SET<Integer>>();
      for (int i = 0; i < words.length; i++)
      {
         String s = words[i];
         if (!st.contains(s))
            st.put(s, new SET<Integer>());
         SET<Integer> pages = st.get(s);
         set.put(i);
      }

      while (!StdIn.isEmpty())
      {
         String query = StdIn.readString();
         SET<Integer> set = st.get(query);
         for (int k : set)
            // print words[k-5] to words[k+5]
      }
   }
}
```

read text and
build index

process queries
and print
concordances

# Matrix-vector multiplication (standard implementation)

$$
\begin{array}{ccc}
\text{a[][]} & \text{x[]} & \text{b[]} \\
\begin{bmatrix}
0 & .90 & 0 & 0 & 0 \\
0 & 0 & .36 & .36 & .18 \\
0 & 0 & 0 & .90 & 0 \\
.90 & 0 & 0 & 0 & 0 \\
.47 & 0 & .47 & 0 & 0
\end{bmatrix}
\begin{bmatrix}
.05 \\
.04 \\
.36 \\
.37 \\
.19
\end{bmatrix}
=
\begin{bmatrix}
.036 \\
.297 \\
.333 \\
.045 \\
.1927
\end{bmatrix}
\end{array}
$$

```
...
double[][] a = new double[N][N];
double[] x = new double[N];
double[] b = new double[N];
...
// initialize a[][] and x[]
...
for (int i = 0; i < N; i++)          nested loops
{                                    N² running time
    sum = 0.0;
    for (int j = 0; j < N; j++)
        sum += a[i][j]*x[j];
    b[i] = sum;
}
```

## Sparse matrix-vector multiplication

Problem.  Sparse matrix-vector multiplication.

Assumptions.  Matrix dimension is 10,000; average nonzeros per row ~ 10.



A    *    x    =    b

## Vector representations

1D array (standard) representation.

- Constant time access to elements.
- Space proportional to N.

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|
| 0 | .36 | 0 | 0 | 0 | .36 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | .18 | 0 | 0 | 0 | 0 | 0 |

Symbol table representation.

- key = index, value = entry
- Efficient iterator.
- Space proportional to number of nonzeros.

## Sparse vector data type

```
public class SparseVector
{
    private HashST<Integer, Double> v;                    ← HashST because order not important

    public SparseVector()
    {   v = new HashST<Integer, Double>();   }            ← empty ST represents all 0s vector

    public void put(int i, double x)
    {   v.put(i, x);   }                                  ← a[i] = value

    public double get(int i)
    {
        if (!v.contains(i)) return 0.0;
        else return v.get(i);                             ← return a[i]
    }

    public Iterable<Integer> indices()
    {   return v.keys();   }

    public double dot(double[] that)
    {
        double sum = 0.0;                                 ← dot product is constant
        for (int i : indices())                             time for sparse vectors
            sum += that[i]*this.get(i);
        return sum;
    }
}
```

24

# Matrix representations

**2D array (standard) representation:** Each row of matrix is an array.

- Constant time access to elements.
- Space proportional to $N^2$.

**Sparse representation:** Each row of matrix is a sparse vector.

- Efficient access to elements.
- Space proportional to number of nonzeros (plus N).

## Sparse matrix-vector multiplication

$$
\begin{array}{ccc}
\texttt{a[][]} & \texttt{x[]} & \texttt{b[]} \\[4pt]
\begin{bmatrix}
0 & .90 & 0 & 0 & 0 \\
0 & 0 & .36 & .36 & .18 \\
0 & 0 & 0 & .90 & 0 \\
.90 & 0 & 0 & 0 & 0 \\
.47 & 0 & .47 & 0 & 0
\end{bmatrix}
&
\begin{bmatrix}
.05 \\
.04 \\
.36 \\
.37 \\
.19
\end{bmatrix}
=
&
\begin{bmatrix}
.036 \\
.297 \\
.333 \\
.045 \\
.1927
\end{bmatrix}
\end{array}
$$

```
..
SparseVector[] a;
a = new SparseVector[N];
double[] x = new double[N];
double[] b = new double[N];
...
// Initialize a[] and x[]
...
for (int i = 0; i < N; i++)
    b[i] = a[i].dot(x);
```

one loop
linear running time
for sparse matrix

Searching challenge 2A:

Problem.  IP lookups in a web monitoring device.

Assumption A.  Billions of lookups, millions of distinct addresses.

Which searching method to use?
1)  Sequential search in a linked list.
2)  Binary search in an ordered array.
3)  Need better method, all too slow.
4)  Doesn't matter much, all fast enough.

Problem. IP lookups in a web monitoring device.

Assumption A. Billions of lookups, millions of distinct addresses.

Which searching method to use?

1) Sequential search in a linked list.

2) Binary search in an ordered array.

✓ 3) Need better method, all too slow.

4) Doesn't matter much, all fast enough.

total cost of insertions is $c*1000000^2$ = $c*1{,}000{,}000{,}000{,}000$ (way too much)

## Searching challenge 2B

Problem.  IP lookups in a web monitoring device.

Assumption B.  Billions of lookups, thousands of distinct addresses.

Which searching method to use?

1) Sequential search in a linked list.
2) Binary search in an ordered array.
3) Need better method, all too slow.
4) Doesn't matter much, all fast enough.

**Problem.** IP lookups in a web monitoring device.

**Assumption B.** Billions of lookups, thousands of distinct addresses.

Which searching method to use?

1) Sequential search in a linked list.

✓ 2) Binary search in an ordered array. ⟵

3) Need better method, all too slow.

4) Doesn't matter much, all fast enough.

total cost of insertions is
$c_1*1000^2 = c_1*1000000$
and dominated by $c_2*1000000000$
cost of lookups

Problem.  Spell checking for a book.

Assumptions.  Dictionary has 25,000 words; book has 100,000+ words.

Which searching method to use?

1)  Sequential search in a linked list.
2)  Binary search in an ordered array.
3)  Need better method, all too slow.
4)  Doesn't matter much, all fast enough.

Problem.  Spell checking for a book.

Assumptions.  Dictionary has 25,000 words; book has 100,000+ words.

Which searching method to use?

1) Sequential search in a linked list.

✓ 2) Binary search in an ordered array. ← easy to presort dictionary total cost of lookups is optimal $c_2*1,500,000$

3) Need better method, all too slow.

4) Doesn't matter much, all fast enough.

Problem.  Maintain symbol table of song names for an iPod.

Assumption A.  Hundreds of songs.

Which searching method to use?

1) Sequential search in a linked list.
2) Binary search in an ordered array.
3) Need better method, all too slow.
4) Doesn't matter much, all fast enough.

Searching challenge 1A

Problem.  Maintain symbol table of song names for an iPod.

Assumption A.  Hundreds of songs.

Which searching method to use?
1)  Sequential search in a linked list.
2)  Binary search in an ordered array.
3)  Need better method, all too slow.
✓ 4)  Doesn't matter much, all fast enough.  ⟵  $100^2 = 10{,}000$

Searching challenge 1B

Problem.  Maintain symbol table of song names for an iPod.

Assumption B.  Thousands of songs.

Which searching method to use?
1)  Sequential search in a linked list.
2)  Binary search in an ordered array.
3)  Need better method, all too slow.
4)  Doesn't matter much, all fast enough.

## Searching challenge 1B

Problem.  Maintain symbol table of song names for an iPod.

Assumption B.  Thousands of songs.

Which searching method to use?

1) Sequential search in a linked list.
2) Binary search in an ordered array.  ← maybe, but $1000^2 = 1,000,000$ so user might wait for complete rebuild of index
✓ 3) Need better method, all too slow.
4) Doesn't matter much, all fast enough.

## Searching challenge 3

Problem. Frequency counts in "Tale of Two Cities."

Assumptions. Book has 135,000+ words; about 10,000 distinct words.

Which searching method to use?

1) Sequential search in a linked list.
2) Binary search in an ordered array.
3) Need better method, all too slow.
4) Doesn't matter much, all fast enough.

## Searching challenge 3

**Problem.** Frequency counts in "Tale of Two Cities."

**Assumptions.** Book has 135,000+ words; about 10,000 distinct words.

**Which searching method to use?**

1) Sequential search in a linked list.
2) Binary search in an ordered array.
✓ 3) Need better method, all too slow.
4) Doesn't matter much, all fast enough.

total cost of searches:
$c_2*1,350,000,000$

maybe, but total cost of insertions is $c_1*100,000,000$

Searching challenge 3 (revisited):

Problem. Frequency counts in "Tale of Two Cities"

Assumptions. Book has 135,000+ words; about 10,000 distinct words.

Which searching method to use?

1) Sequential search in a linked list.

2) Binary search in an ordered array.

3) Need better method, all too slow.

4) Doesn't matter much, all fast enough.

✓ 5) BSTs.

insertion cost < 10000 * 1.38 * lg 10000 < .2 million
lookup cost < 135000 * 1.38 * lg 10000 < 2.5 million

Problem.  Index for a PC or the web.

Assumptions.  1 billion++ words to index.

Which searching method to use?

- Hashing
- Red-black-trees
- Doesn't matter much.

**Problem.** Index for a PC or the web.

**Assumptions.** 1 billion++ words to index.

**Which searching method to use?**

✓ • Hashing
   • Red-black-trees  ←  too much space
   • Doesn't matter much.

**Solution.** Symbol table with:

• Key = query string.
• Value = set of pointers to files.

sort the (relatively few) search hits

| Spotlight | searching challenge   ⊗ |
|---|---|
| | Show All (200) |
| **Top Hit** | 10Hashing |
| Documents | mobydick.txt |
| | movies.txt |
| | Papers/Abstracts |
| | score.card.txt |
| | Requests |
| Mail Messages | Re: Draft of lecture on symb… |
| | SODA 07 Final Accepts |
| | SODA 07 Summary |
| | Got–it |
| | No Subject |
| PDF Documents | 08BinarySearchTrees.pdf |
| | 07SymbolTables.pdf |
| | 07SymbolTables.pdf |
| | 06PriorityQueues.pdf |
| | 06PriorityQueues.pdf |
| Presentations | 10Hashing |
| | 07SymbolTables |
| | 06PriorityQueues |

Problem. Index for an e-book.

Assumptions. Book has 100,000+ words.

Which searching method to use?

1. Hashing
2. Red-black-tree
3. Doesn't matter much.

**Problem.** Index for an e-book.

**Assumptions.** Book has 100,000+ words.

**Which searching method to use?**

1. Hashing
✓ 2. Red-black-tree ← need ordered iteration
3. Doesn't matter much.

**Solution.** Symbol table with:

- Key = index term.
- Value = ordered set of pages on which term appears.