

# Summary & Review

Dr. 何明晰, He Mingxin, Max

program06 @ yeah.net

Email Subject: CS203B + *ID* + *Name*: *TOPIC*

Sakai: CS203B Fall 2022

数据结构与算法分析B

Data Structures and Algorithm Analysis

# Key Terms

## Algorithm

- A **problem-solving method** suitable for implementation as a computer program

## Data structures

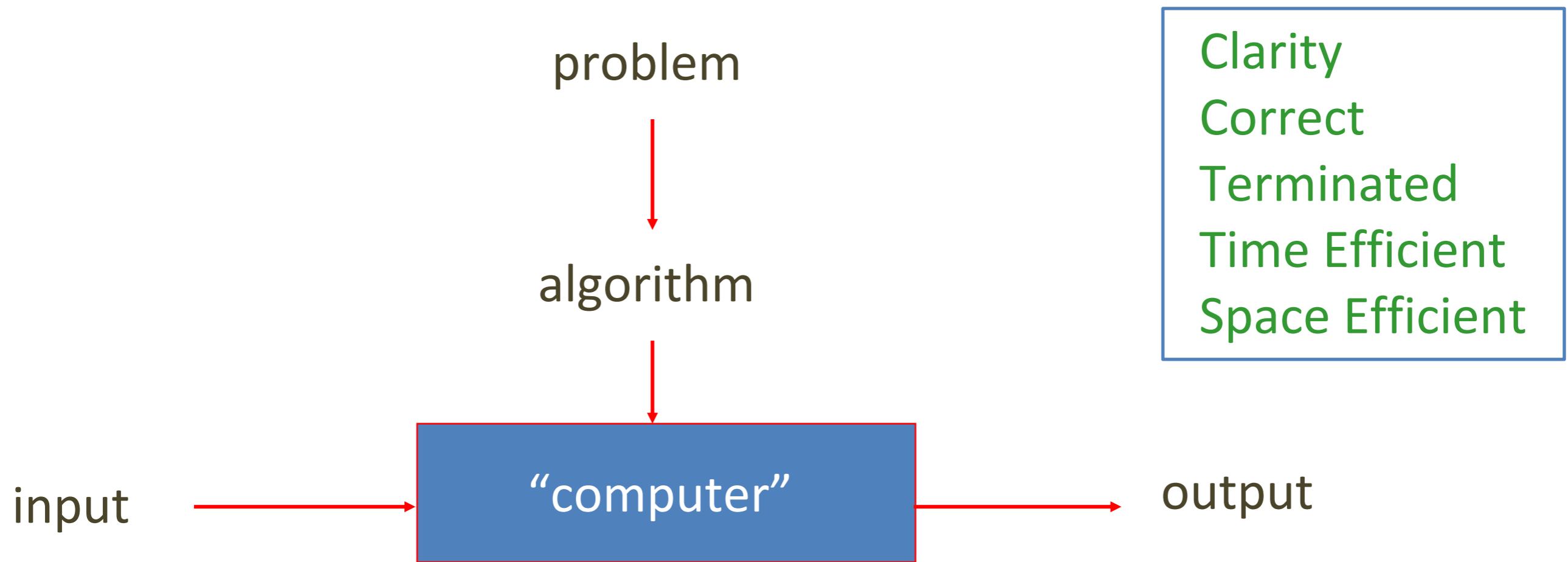
- Objects created to organize data used in computation
- **A way to store and organize data in order to facilitate access and modifications**

Data structure exist as the by-product or end product of algorithms

- Understanding data structure is essential to understanding algorithms and hence to problem-solving
- Simple algorithms can give rise to complicated data-structures
- Complicated algorithms can use simple data structures

# What is an Algorithm?

An *algorithm* is a sequence of unambiguous instructions for solving a computation problem, i.e., for obtaining a required output for any legitimate input in a finite amount of time.



# Why study Data Structures and algorithms?

## Using a computer?

- Solve computational problems?
- Want it to go faster?
- Ability to process more data?

## Technology vs. Performance/cost factor

- Technology can improve things by a constant factor
- Good algorithm design can do much better and may be cheaper
- Supercomputer cannot rescue a bad algorithm

## Data structures and algorithms as a field of study

- Old enough to have basics known
- New discoveries
- Burgeoning application areas
- Philosophical implications?

# Data Structures covered in the course

Data Structures	Advantages	Disadvantages
Array	Quick insertion, very fast access if index known	Slow search (fast after sorting), slow deletion, fixed size.
Ordered array	Quicker search than unsorted array.	Slow insertion and deletion, fixed size.
Stack	Provides last-in, first-out access.	Slow access to other items.
Queue	Provides first-in, first-out access.	Slow access to other items.
Linked list	Quick insertion, quick deletion.	Slow search.
Binary tree	Quick search, insertion, deletion	Deletion algorithm is complex.
Graph	Models real-world situations.	Some algorithms are slow and complex.
Hashing	Quick search, insertion, and deletion	Collision

# Topics Covered:

[15UnionFind](#)

[14AnalysisOfAlgorithms](#)

[13StacksAndQueues](#)

[21ElementarySorts](#)

[22Mergesort](#)

[23Quicksort](#)

[24PriorityQueues](#)

[A31ElementarySymbolTables](#)

[B32BinarySearchTrees](#)

[33BalancedSearchTrees](#)

[41UndirectedGraphs](#)

[42DirectedGraphs](#)

[43MinimumSpanningTrees](#)

[51StringSorts](#)

[53SubstringSearch](#)

# Course Evaluation:

- Class participation/quiz 10%
- Lab 20%
- Course project 20%
- Mid term exam 20%
- Final exam 30%

# Reasons to analyze algorithms

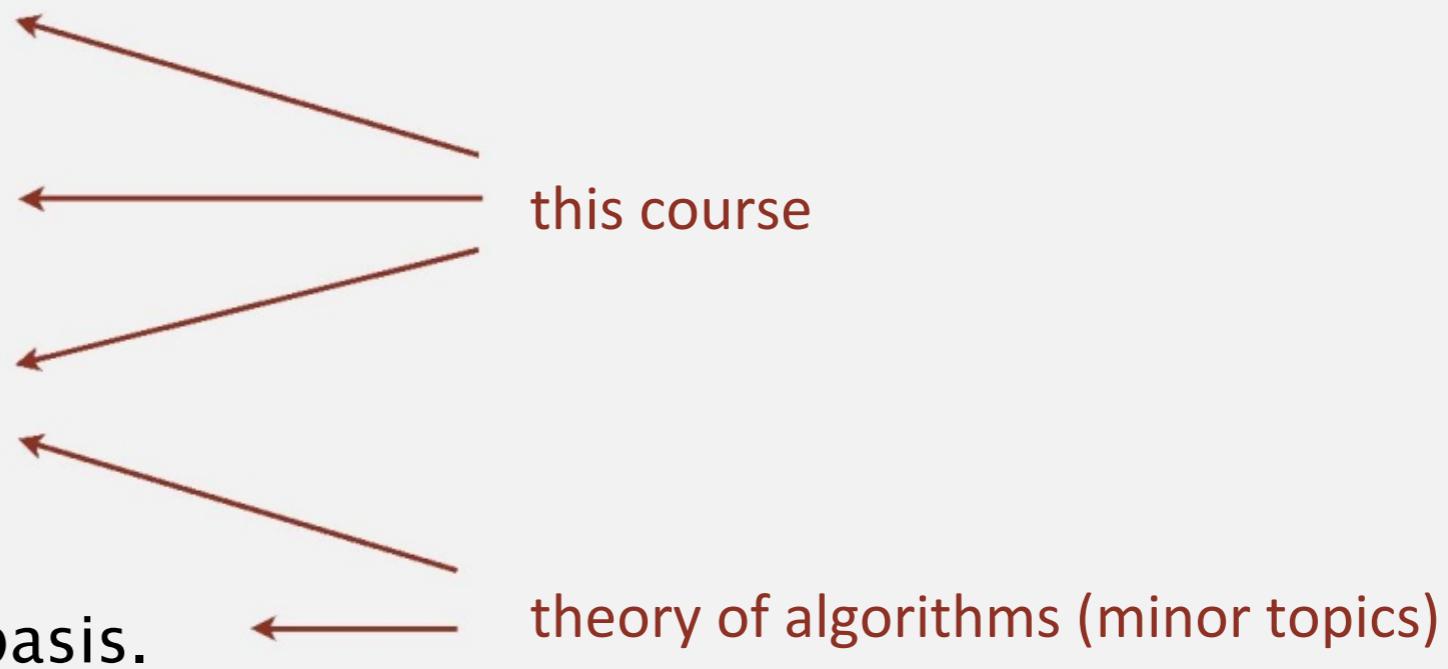
---

Predict performance.

Compare algorithms.

Provide guarantees.

Understand theoretical basis.



Primary practical reason: avoid performance bugs.



client gets poor performance because programmer  
did not understand performance characteristics



# Scientific method applied to analysis of algorithms

---

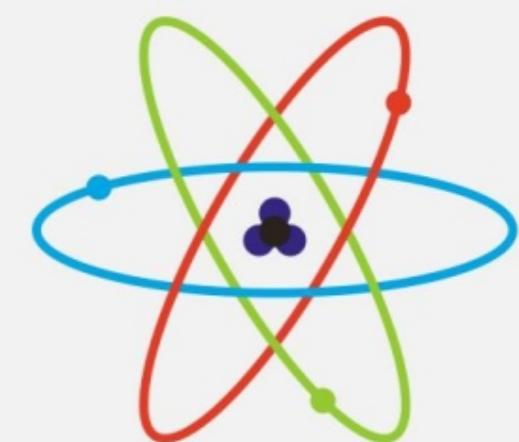
A framework for predicting performance and comparing algorithms.

## Scientific method.

- **Observe** some feature of the natural world.
- **Hypothesize** a model that is consistent with the observations.
- **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**.



Feature of the natural world. Computer itself.

# Experimental algorithmics

---

## System independent effects.

- Algorithm.
  - Input data.
- 
- determines exponent  
in power law

## System dependent effects.

- Hardware: CPU, memory, cache, ...
  - Software: compiler, interpreter, garbage collector, ...
  - System: operating system, network, other apps, ...
- 
- determines constant  
in power law

Bad news. Difficult to get precise measurements.

Good news. Much easier and cheaper than other sciences.



e.g., can run huge number of experiments

## Simplification 1: cost model

Cost model. Use some basic operation as a proxy for running time.

```
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	$\frac{1}{2} (N + 1) (N + 2)$
equal to compare	$\frac{1}{2} N (N - 1)$
<b>array access</b>	$N (N - 1)$
increment	$\frac{1}{2} N (N - 1)$ to $N (N - 1)$

$$\begin{aligned}0 + 1 + 2 + \dots + (N - 1) &= \frac{1}{2} N (N - 1) \\&= \binom{N}{2}\end{aligned}$$

cost model = array accesses  
(we assume compiler/JVM do not optimize any array accesses away!)

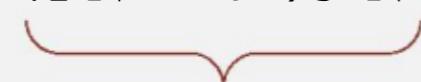
## Simplification 2: tilde notation

- 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

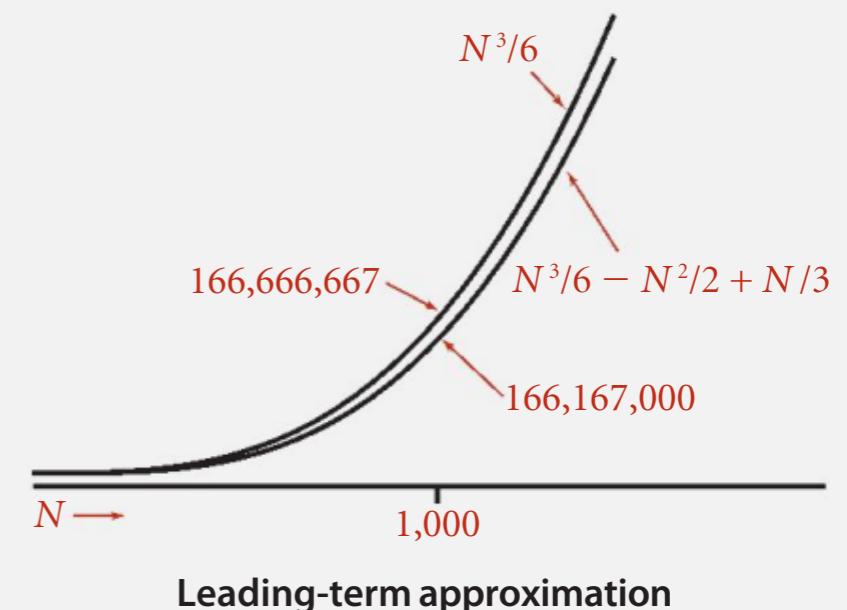
Ex 1.  $\frac{1}{6}N^3 + 20N + 16 \sim \frac{1}{6}N^3$

Ex 2.  $\frac{1}{6}N^3 + 100N^{4/3} + 56 \sim \frac{1}{6}N^3$

Ex 3.  $\frac{1}{6}N^3 - \frac{1}{2}N^2 + \frac{1}{3}N \sim \frac{1}{6}N^3$

 discard lower-order terms

(e.g.,  $N = 1000$ : 166.67 million vs. 166.17 million)



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

## Simplification 2: 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

operation	frequency	tilde notation
variable declaration	$N + 2$	$\sim N$
assignment statement	$N + 2$	$\sim N$
less than compare	$\frac{1}{2} (N + 1)(N + 2)$	$\sim \frac{1}{2} N^2$
equal to compare	$\frac{1}{2} N(N - 1)$	$\sim \frac{1}{2} N^2$
array access	$N(N - 1)$	$\sim N^2$
increment	$\frac{1}{2} N(N - 1)$ to $N(N - 1)$	$\sim \frac{1}{2} N^2$ to $\sim N^2$

## Diversion: estimating a discrete sum

---

Q. How to estimate a discrete sum?

A1. Take a discrete mathematics course.

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

Ex 1.  $1 + 2 + \dots + N$ .

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

Ex 2.  $1^k + 2^k + \dots + N^k$ .

$$\sum_{i=1}^N i^k \sim \int_{x=1}^N x^k dx \sim \frac{1}{k+1} N^{k+1}$$

Ex 3.  $1 + 1/2 + 1/3 + \dots + 1/N$ .

$$\sum_{i=1}^N \frac{1}{i} \sim \int_{x=1}^N \frac{1}{x} dx = \ln N$$

Ex 4. 3-sum triple loop.

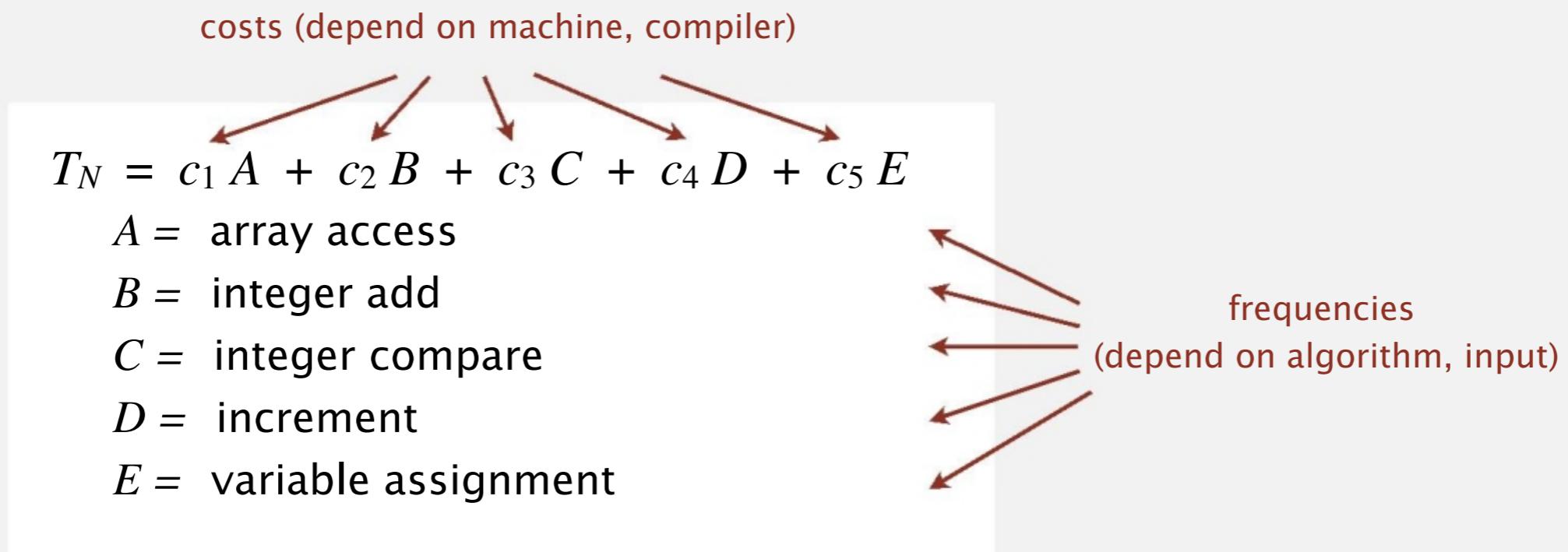
$$\sum_{i=1}^N \sum_{j=i}^N \sum_{k=j}^N 1 \sim \int_{x=1}^N \int_{y=x}^N \int_{z=y}^N dz dy dx \sim \frac{1}{6} N^3$$

# 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.



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

## Common order-of-growth classifications

---

**Definition.** If  $f(N) \sim c g(N)$  for some constant  $c > 0$ , then the order of growth of  $f(N)$  is  $g(N)$ .

- Ignores leading coefficient.
- Ignores lower-order terms.

**Ex.** The order of growth of the running time of this code is  $N^3$ .

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

**Typical usage.** With running times.



where leading coefficient depends on  
machine, compiler, JVM, ...

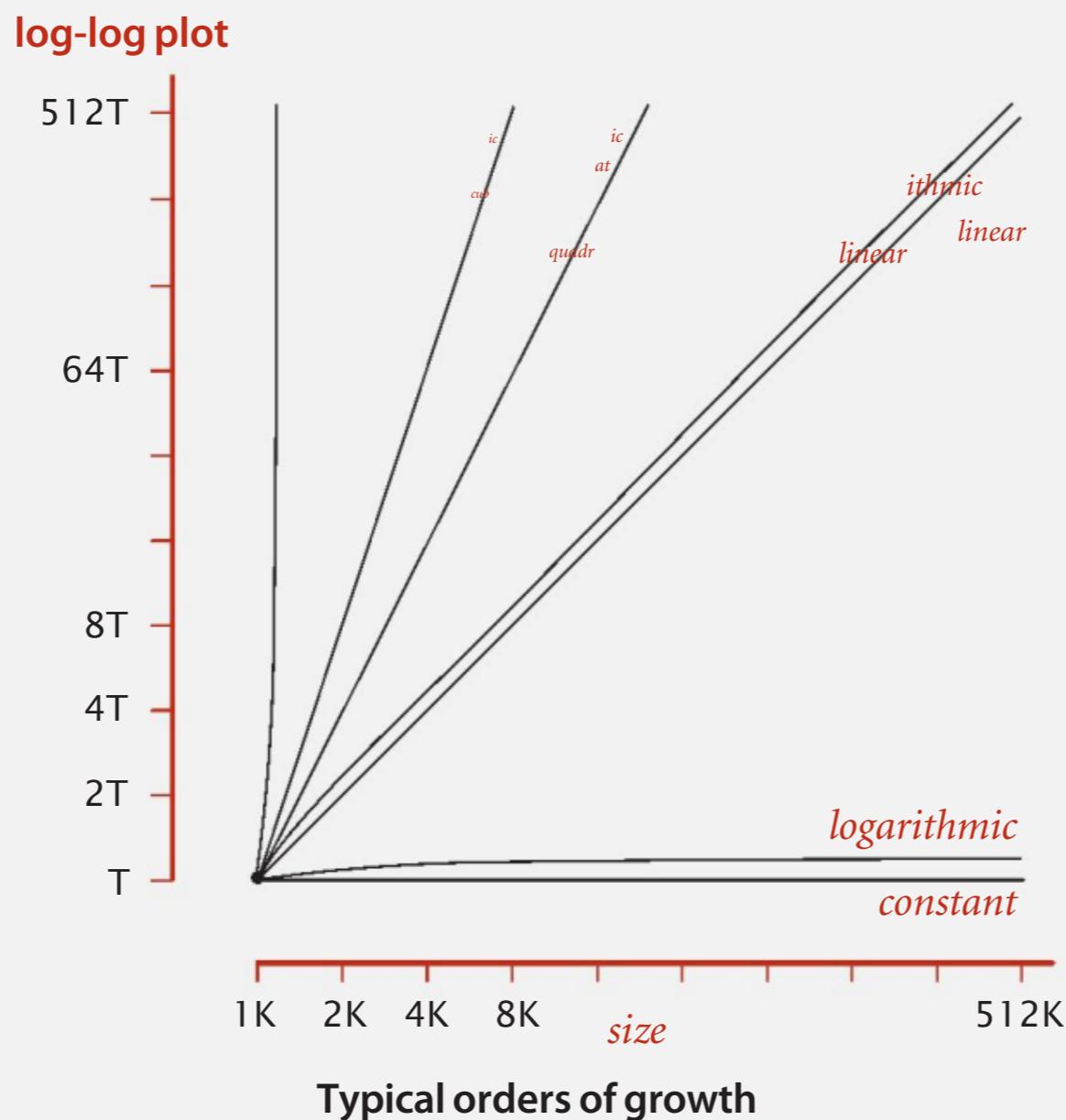
# Common order-of-growth classifications

---

Good news. The set of functions

$1, \log N, N, N \log N, N^2, N^3,$  and  $2^N$

suffices to describe the order of growth of most common algorithms.



# Common order-of-growth classifications

---

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

## Binary search: mathematical analysis

---

**Proposition.** Binary search uses at most  $1 + \lg N$  key compares to search in a sorted array of size  $N$ .

**Def.**  $T(N) = \# \text{ key compares to binary search a sorted subarray of size } \leq N$ .

**Binary search recurrence.**  $T(N) \leq T(N/2) + 1$  for  $N > 1$ , with  $T(1) = 1$ .

$\uparrow$        $\uparrow$   
left or right half      possible to implement with one  
(floored division)      2-way compare (instead of 3-way)

**Pf sketch.** [assume  $N$  is a power of 2]

$$\begin{aligned} T(N) &\leq T(N/2) + 1 && [\text{given}] \\ &\leq T(N/4) + 1 + 1 && [\text{apply recurrence to first term}] \\ &\leq T(N/8) + 1 + 1 + 1 && [\text{apply recurrence to first term}] \\ &\vdots \\ &\leq T(N/N) + 1 + 1 + \dots + 1 && [\text{stop applying, } T(1) = 1] \\ &= 1 + \lg N \end{aligned}$$

# 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 a guarantee for all inputs.

**Average case.** Expected cost for random input.

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

this course

**Ex 1.** Array accesses for brute-force 3-SUM.

Best:  $\sim \frac{1}{2} N^3$

Average:  $\sim \frac{1}{2} N^3$

Worst:  $\sim \frac{1}{2} N^3$

**Ex 2.** Comparisons for binary search.

Best:  $\sim 1$

Average:  $\sim \lg N$

Worst:  $\sim \lg N$

# Theory of algorithms

---

## Goals.

- Establish “difficulty/complexity” of a problem.
- Develop “optimal” algorithms.

## Approach.

- Suppress details in analysis: analyze “to within a constant factor.”
- Eliminate variability in input model: focus on the worst case.

**Upper bound.** Performance guarantee of algorithm for any input.

**Lower bound.** Proof that no algorithm can do better.

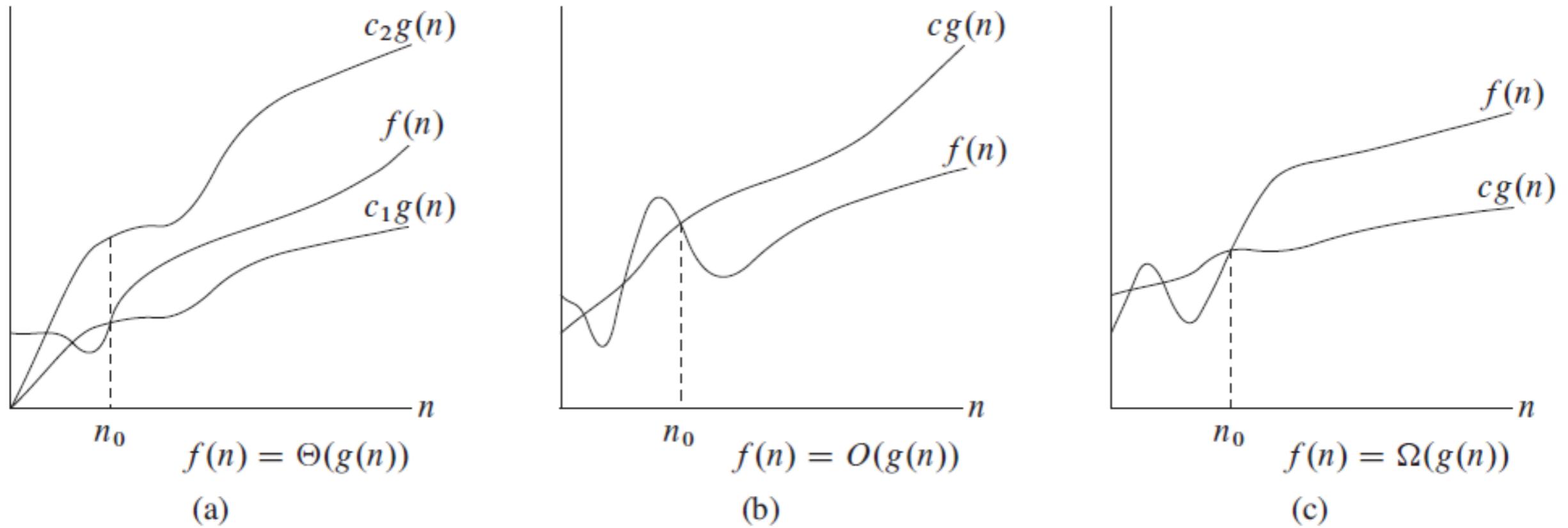
**Optimal algorithm.** Lower bound = upper bound (to within a constant factor).

# Commonly-used notations in the theory of algorithms

---

notation	provides	example	shorthand for	used to
<b>Big Theta</b>	asymptotic order of growth	$\Theta(N^2)$	$\frac{1}{2} N^2$ $10 N^2$ $5 N^2 + 22 N \log N + 3N$ ⋮	classify algorithms
<b>Big Oh</b>	$\Theta(N^2)$ and smaller	$O(N^2)$	$10 N^2$ $100 N$ $22 N \log N + 3 N$ ⋮	develop upper bounds
<b>Big Omega</b>	$\Theta(N^2)$ and larger	$\Omega(N^2)$	$\frac{1}{2} N^2$ $N^5$ $N^3 + 22 N \log N + 3 N$ ⋮	develop lower bounds

# Asymptotic notation (渐近记号)



**Figure 3.1** Graphic examples of the  $\Theta$ ,  $O$ , and  $\Omega$  notations. In each part, the value of  $n_0$  shown is the minimum possible value; any greater value would also work. **(a)**  $\Theta$ -notation bounds a function to within constant factors. We write  $f(n) = \Theta(g(n))$  if there exist positive constants  $n_0$ ,  $c_1$ , and  $c_2$  such that at and to the right of  $n_0$ , the value of  $f(n)$  always lies between  $c_1 g(n)$  and  $c_2 g(n)$  inclusive. **(b)**  $O$ -notation gives an upper bound for a function to within a constant factor. We write  $f(n) = O(g(n))$  if there are positive constants  $n_0$  and  $c$  such that at and to the right of  $n_0$ , the value of  $f(n)$  always lies on or below  $c g(n)$ . **(c)**  $\Omega$ -notation gives a lower bound for a function to within a constant factor. We write  $f(n) = \Omega(g(n))$  if there are positive constants  $n_0$  and  $c$  such that at and to the right of  $n_0$ , the value of  $f(n)$  always lies on or above  $c g(n)$ .

# Theory of algorithms: example 1

---

## Goals.

- Establish “difficulty” of a problem and develop “optimal” algorithms.
- Ex. 1-SUM = “*Is there a 0 in the array?* ”

## Upper bound.

A specific algorithm.

- Ex. Brute-force algorithm for 1-SUM: Look at every array entry.
- Running time of the optimal algorithm for 1-SUM is  $O(N)$ .

## Lower bound.

Proof that no algorithm can do better.

- Ex. Have to examine all  $N$  entries (any unexamined one might be 0).
- Running time of the optimal algorithm for 1-SUM is  $\Omega(N)$ .

## Optimal algorithm.

- Lower bound equals upper bound (to within a constant factor).
- Ex. Brute-force algorithm for 1-SUM is optimal: its running time is  $\Theta(N)$ .

# Commonly-used notations in the theory of algorithms

notation	provides	example	shorthand for	used to
<b>Tilde</b>	leading term	$\sim 10 N^2$	$10 N^2$ $10 N^2 + 22 N \log N$ $10 N^2 + 2 N + 37$	provide approximate model
<b>Big Theta</b>	asymptotic order of growth	$\Theta(N^2)$	$\frac{1}{2} N^2$ $10 N^2$ $5 N^2 + 22 N \log N + 3N$	classify algorithms
<b>Big Oh</b>	$\Theta(N^2)$ and smaller	$O(N^2)$	$10 N^2$ $100 N$ $22 N \log N + 3 N$	develop upper bounds
<b>Big Omega</b>	$\Theta(N^2)$ and larger	$\Omega(N^2)$	$\frac{1}{2} N^2$ $N^5$ $N^3 + 22 N \log N + 3 N$	develop lower bounds

**Common mistake.** Interpreting big-Oh as an approximate model.  
**This course.** Focus on approximate models: use Tilde-notation

# Basics

---

Bit. 0 or 1.

NIST

most computer scientists

Byte. 8 bits.



Megabyte (MB). 1 million or  $2^{20}$  bytes.

Gigabyte (GB). 1 billion or  $2^{30}$  bytes.



64-bit machine. We assume a 64-bit machine with 8-byte pointers.

- Can address more memory.
- Pointers use more space.

some JVMs "compress" ordinary object  
pointers to 4 bytes to avoid this cost



# Typical memory usage for primitive types and arrays

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

primitive types

type	bytes
char[]	$2N + 24$
int[]	$4N + 24$
double[]	$8N + 24$

one-dimensional arrays

type	bytes
char[][]	$\sim 2MN$
int[][]	$\sim 4MN$
double[][]	$\sim 8MN$

two-dimensional arrays

# Typical memory usage for objects in Java

---

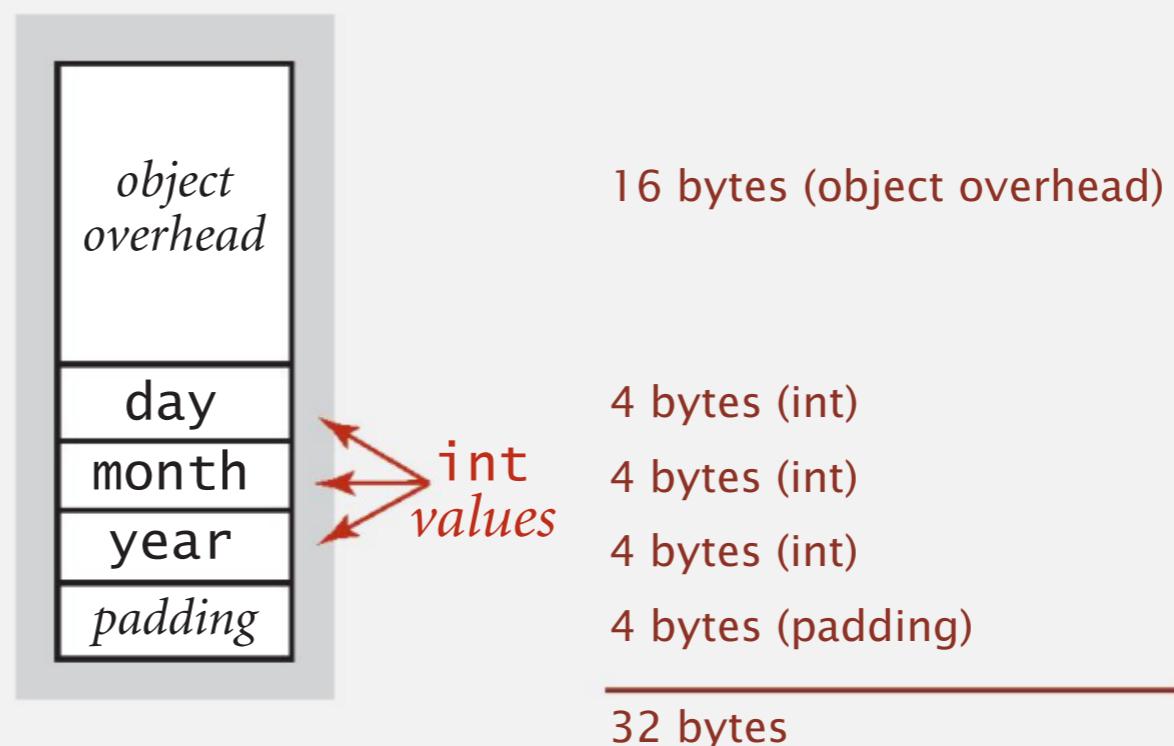
Object overhead. 16 bytes.

Reference. 8 bytes.

Padding. Each object uses a multiple of 8 bytes.

Ex 1. A Date object uses 32 bytes of memory.

```
public class Date {  
  
    private int day;  
    private int month;  
    private int year;  
    ...  
}
```



32 bytes

# Typical memory usage summary

---

## Total memory usage for a data type value:

- Primitive type: 4 bytes for int, 8 bytes for double, ...
- Object reference: 8 bytes.
- Array: 24 bytes + memory for each array entry.
- Object: 16 bytes + memory for each instance variable.
- Padding: round up to multiple of 8 bytes.

+ 8 extra bytes per inner class object  
(for reference to enclosing class)

**Shallow memory usage:** Don't count referenced objects.

**Deep memory usage:** If array entry or instance variable is a reference, count memory (recursively) for referenced object.

## Example

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

Use tilde notation to simplify your answer.

```
public class WeightedQuickUnionUF {  
  
    private int[] id;  
    private int[] sz;  
    private int count;  
  
    public WeightedQuickUnionUF(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;  
    }  
    ...  
}
```

16 bytes  
(object overhead)

8 + (4N + 24) bytes each  
(reference + int[] array)

4 bytes (int)

4 bytes (padding)

---

8N + 88 bytes

A.  $8N + 88 \sim 8N$  bytes.

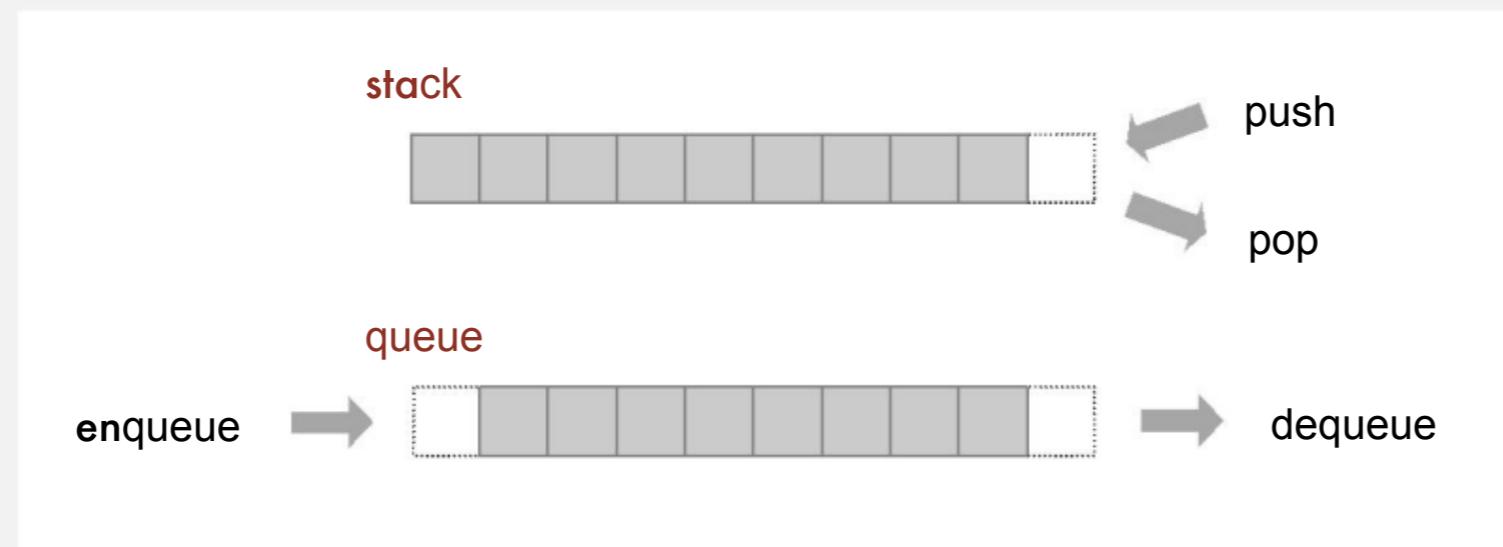
# Algorithm Design Methods: Techniques/Strategies

- Brute Force
  - ▶ Greedy Approach
- Divide and Conquer
  - ▶ Dynamic Programming
- Decrease and Conquer
  - ▶ Iterative improvement
- Transform and Conquer
  - ▶ Backtracking
- Space & Time Tradeoffs
  - ▶ Branch and bound

# Stacks and queues

## Fundamental data types.

- Value: collection of objects.
- Operations: **insert**, **remove**, **iterate**, test if **empty**.
- Intent is clear when we insert.
- Which item do we remove?



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

← LIFO = "last in first out"

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

← FIFO = "first in first out"

# Stack: linked-list implementation in Java

```
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;
    }
}
```

private inner class  
(access modifiers for instance  
variables don't matter)



# Stack: linked-list implementation performance

**Proposition.** Every operation takes constant time in the worst case.

**Proposition.** A stack with N items uses  $\sim 40N$  bytes.

inner class

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



**Remark.** This accounts for the memory for the stack  
(but not the memory for strings themselves, which the client owns).

# Queue dequeue: linked-list implementation

inner class

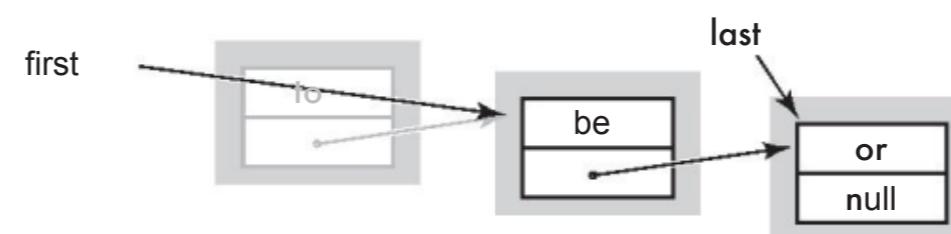
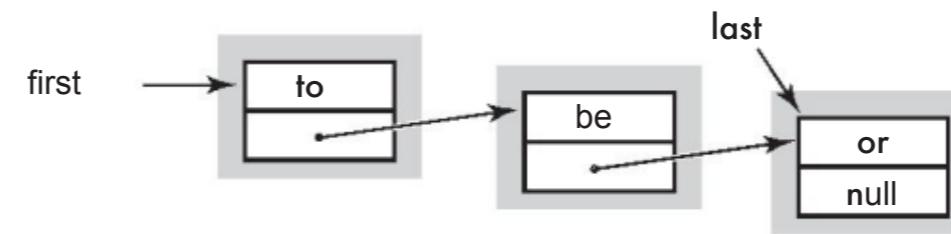
```
private class Node  
{  
    String item;  
    Node next;  
}
```

save item to return

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

delete first node

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



return saved item

```
return item;
```

**Remark.** Identical code to linked-list stack pop().

# Generic stack: array implementation

```
public class FixedCapacityStackOfStrings
{
    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];   }
}
```

the way it should be

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

    public FixedCapacityStack(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];   }
}
```

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

# Generic stack: array implementation

```
public class FixedCapacityStackOfStrings
{
    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];   }
}
```

the way it is

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

    public FixedCapacityStack(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 ugly cast

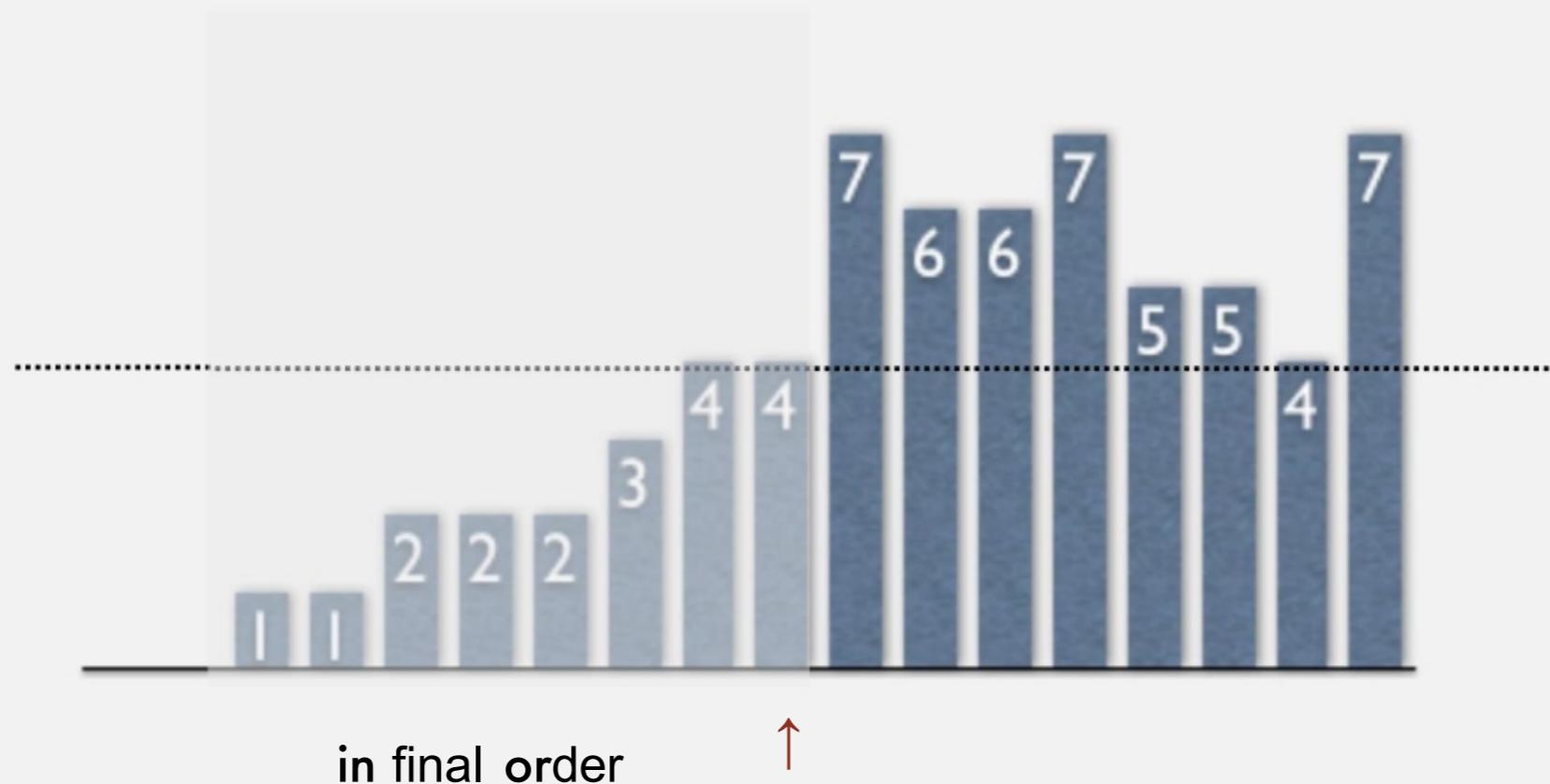
# Selection sort

---

Algorithm.  $\uparrow$  scans from left to right.

Invariants.

- Entries to the left of  $\uparrow$  (including  $\uparrow$ ) fixed and in ascending order.
- No entry to right of  $\uparrow$  is smaller than any entry to the left of  $\uparrow$ .



# Insertion sort

---

Algorithm.  $\uparrow$  scans from left to right.

Invariants.

- Entries to the left of  $\uparrow$  (including  $\uparrow$ ) are in ascending order.
- Entries to the right of  $\uparrow$  have not yet been seen.



# Shellsort overview

---

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

an h-sorted array is h interleaved sorted subsequences



Shellsort. [Shell 1959] **h-sort** array for 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

# Key of Divide and Conquer

- Divide into smaller parts (subproblems)
- Solve the smaller parts **recursively**
- Merge the result of the smaller parts

```
ResultType DandC(Problem p) {  
    if (p is trivial) {  
        solve p directly  
        return the result  
    } else {  
        divide p into p1, p2, ..., pn  
        for (i = 1 to n)  
            ri = DandC(pi)  
        combine r1, r2, ..., rn into r  
        return r  
    }  
}
```

$t_s$

$t_d$

$t_r$

$t_c$

# Mergesort

## Basic plan.

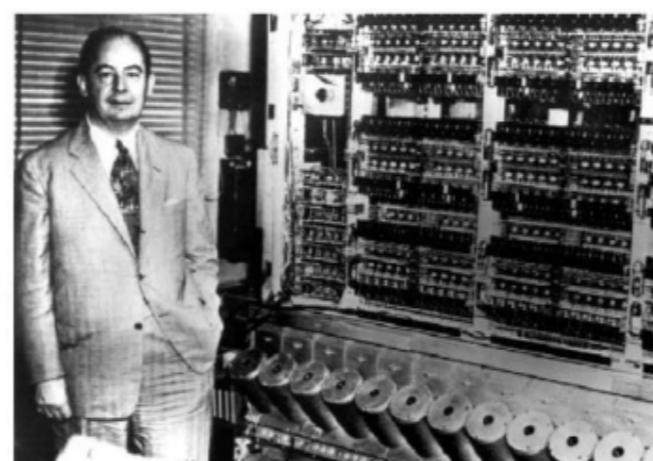
- Divide array into two halves.
- Recursively sort each half.
- Merge two halves.

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



# Quicksort overview

Step 1. Shuffle the array.

Step 2. Partition the array so that, for some  $j$



- Entry  $a[j]$  is in place.
- No larger entry to the left of  $j$ .
- No smaller entry to the right of  $j$ .

Step 3. Sort each subarray recursively.



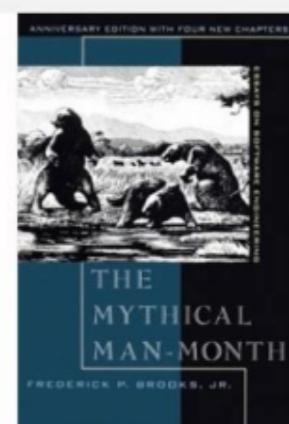
# Collections

---

A **collection** is a data type that stores a group of items.

data type	core operations	data structure
stack	PUSH, POP	linked list, resizing array
queue	ENQUEUE, DEQUEUE	linked list, resizing array
priority queue	INSERT, DELETE-MAX	binary heap
symbol table	PUT, GET, DELETE	binary search tree, hash table
set	ADD, CONTAINS, DELETE	binary search tree, hash table

“Show me your code and conceal your data structures, and I shall continue to be mystified. Show me your data structures, and I won't usually need your code; it'll be obvious.” — Fred Brooks



# Priority queue

---

Collections. Insert and delete items. Which item to delete?

Stack. Remove the item most recently added.

Queue. Remove the item least recently added.

Randomized queue. Remove a random item.

Priority queue. Remove the **largest** (or **smallest**) item.

Generalizes: stack, queue, randomized 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

# Binary heap: representation

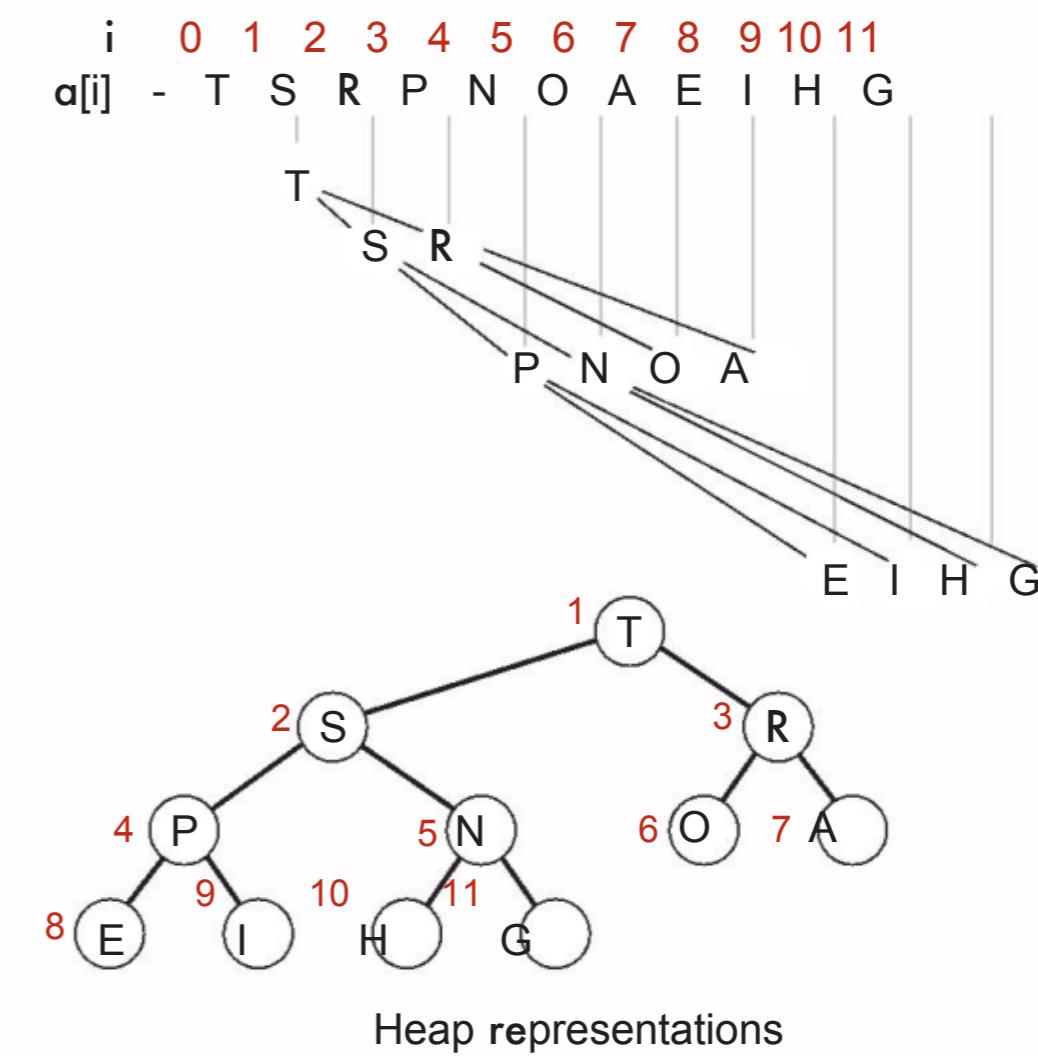
Binary heap. Array representation of a heap-ordered complete binary tree.

Heap-ordered binary tree.

- Keys in nodes.
- Parent's key no smaller than children's keys.

Array representation.

- Indices start at 1.
- Take nodes in **level order**.
- No explicit links needed!

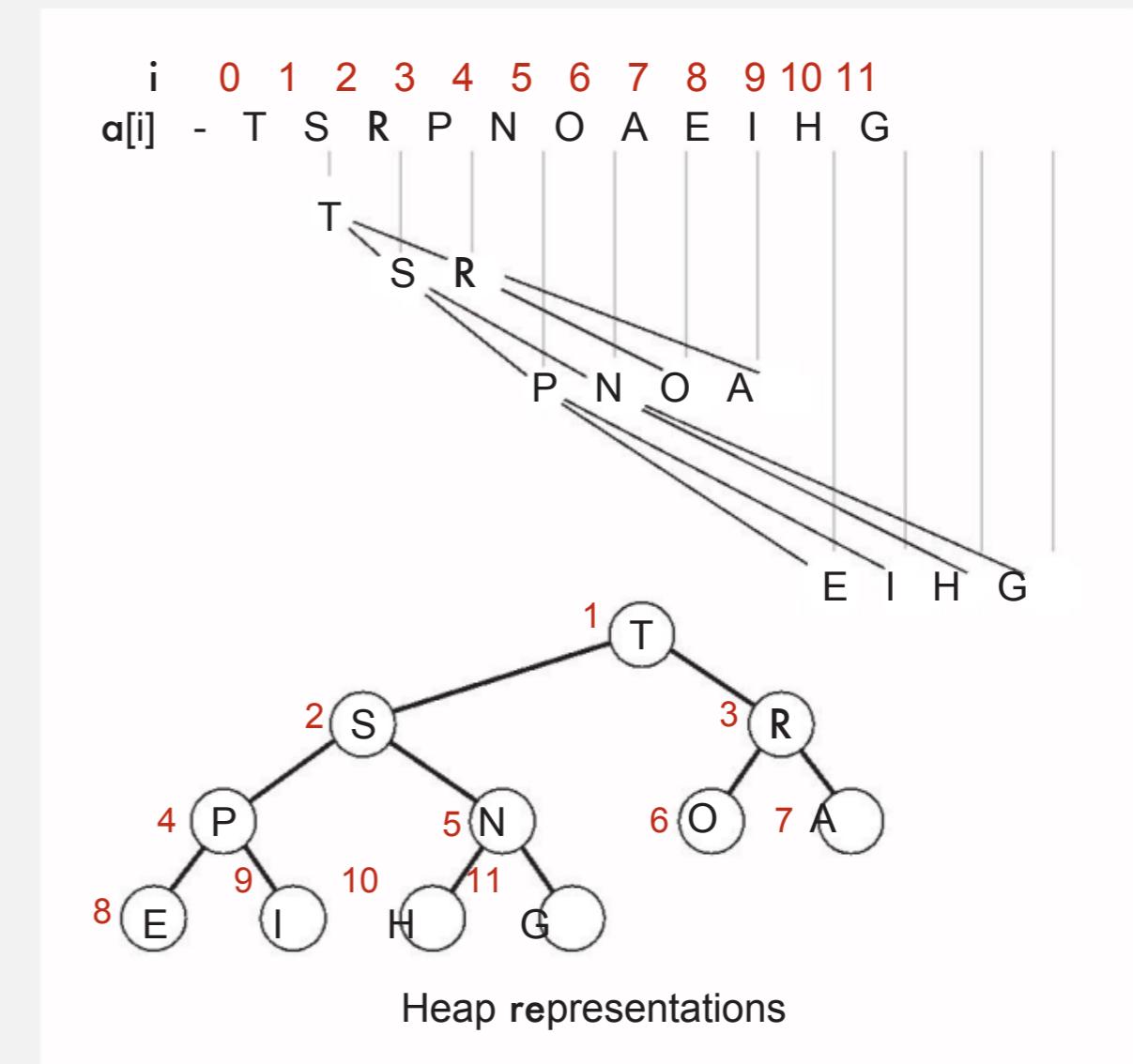


# Binary heap: properties

Proposition. Largest key is  $a[1]$ , which is root of binary tree.

Proposition. 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  $2k+1$ .



# Symbol tables: context

---

Also known as: maps, dictionaries, associative arrays.

Generalizes arrays. Keys need not be between 0 and N – 1.

Language support.

- External libraries: C, VisualBasic, Standard ML, bash, ...
- Built-in libraries: Java, C#, C++, Scala, ...
- Built-in to language: Awk, Perl, PHP, Tcl, JavaScript, Python, Ruby, Lua.

every array is an  
associative array

every object is an  
associative array

table is the only  
primitive data structure

```
hasNiceSyntaxForAssociativeArrays["Python"] = true  
hasNiceSyntaxForAssociativeArrays["Java"]    = false
```

legal Python code

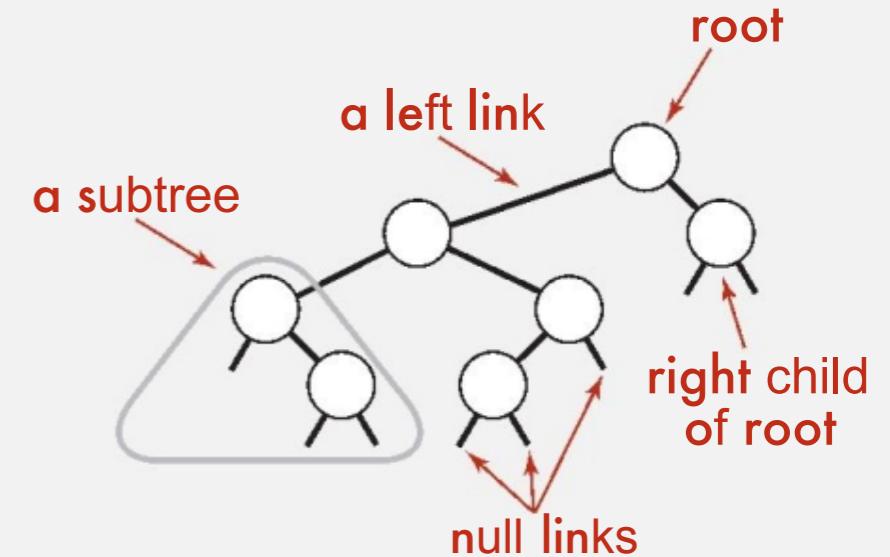
# Binary search trees

---

Definition. A BST is a **binary tree in symmetric order**.

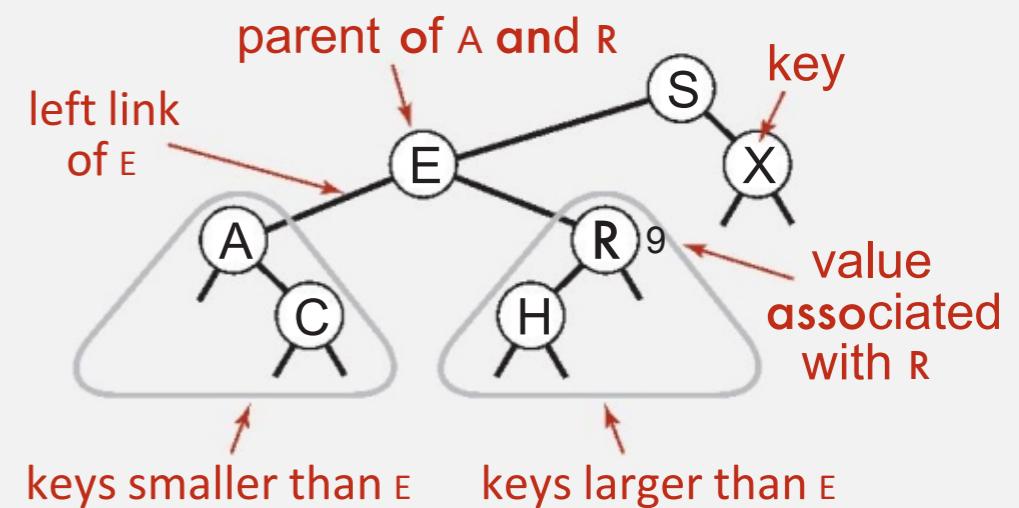
A binary tree is either:

- Empty.
- Two disjoint binary trees (left and right).



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.



## 2-3 tree

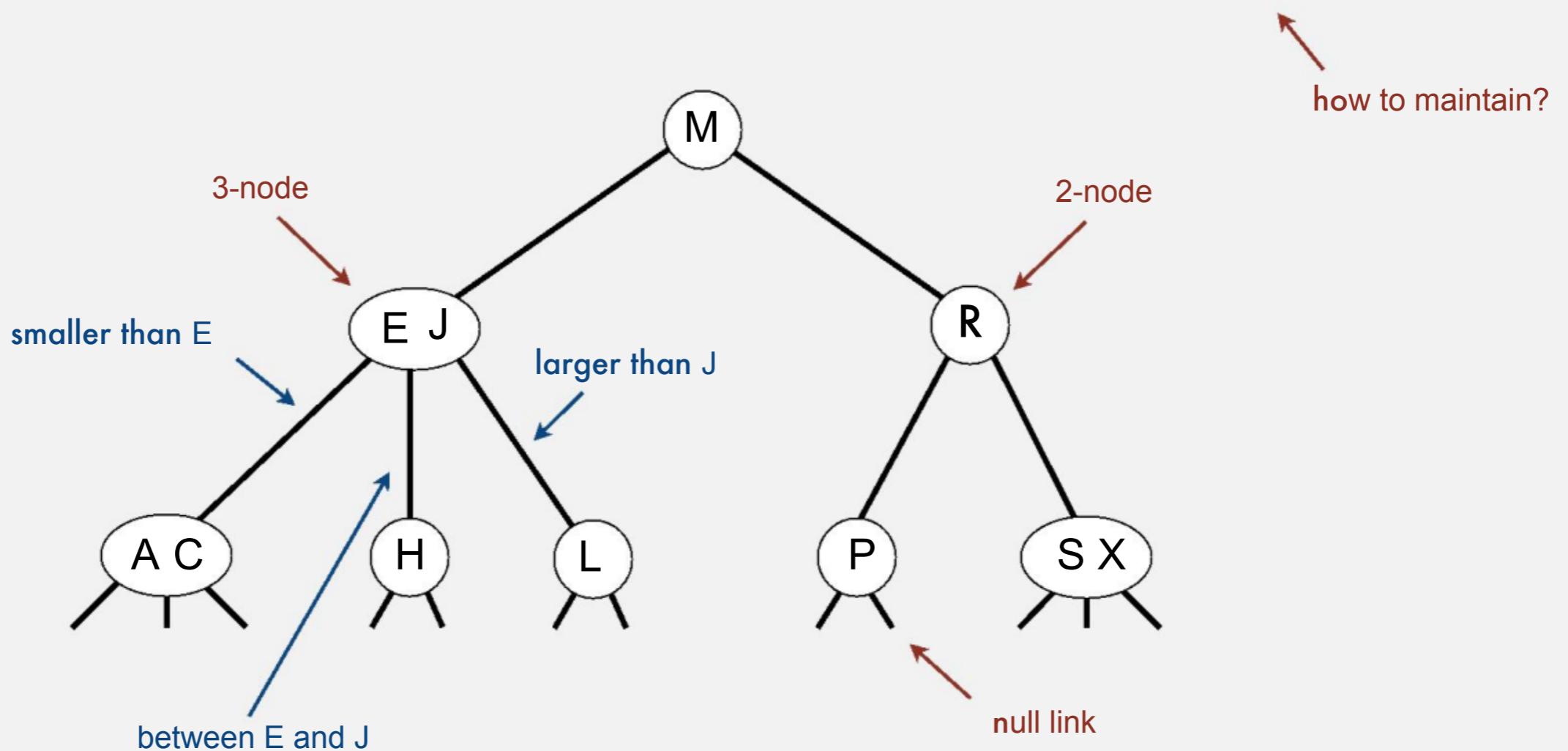
---

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.

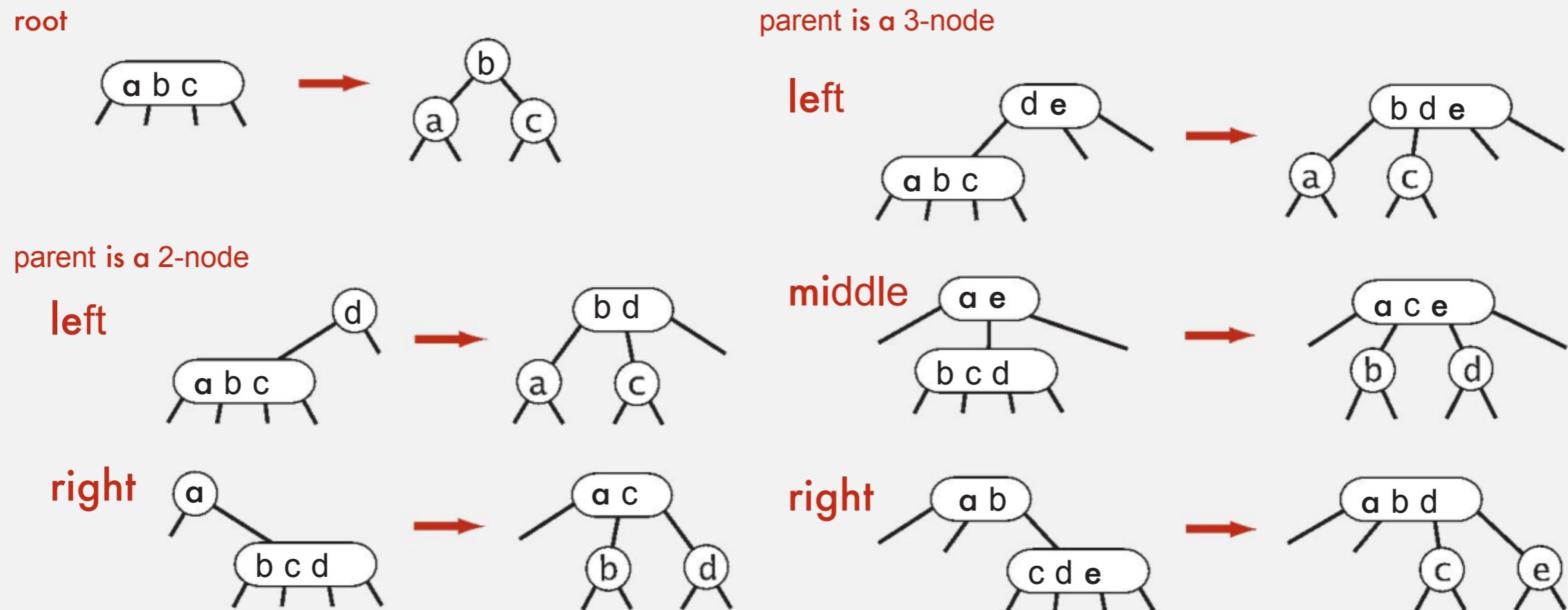
Perfect balance. Every path from root to null link has same length.



# Global properties in a 2-3 tree

Invariants. Maintains symmetric order and perfect balance.

Pf. Each transformation maintains symmetric order and perfect balance.

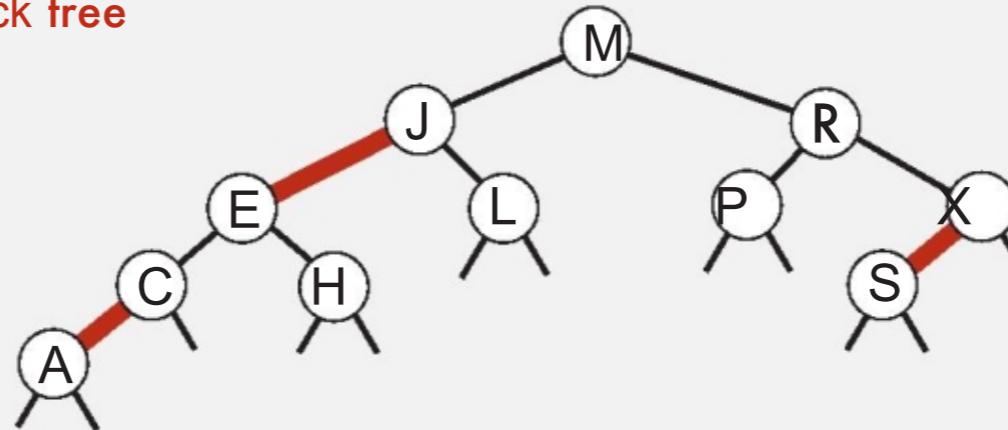


# Left-leaning red-black BSTs: 1-1 correspondence with 2-3 trees

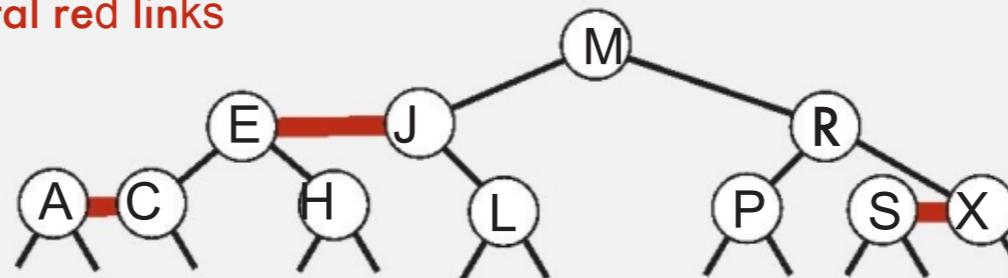
---

Key property. 1–1 correspondence between 2–3 and LLRB.

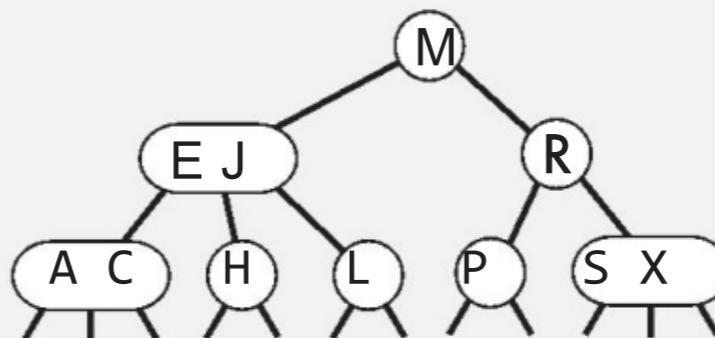
red-black tree



horizontal red links



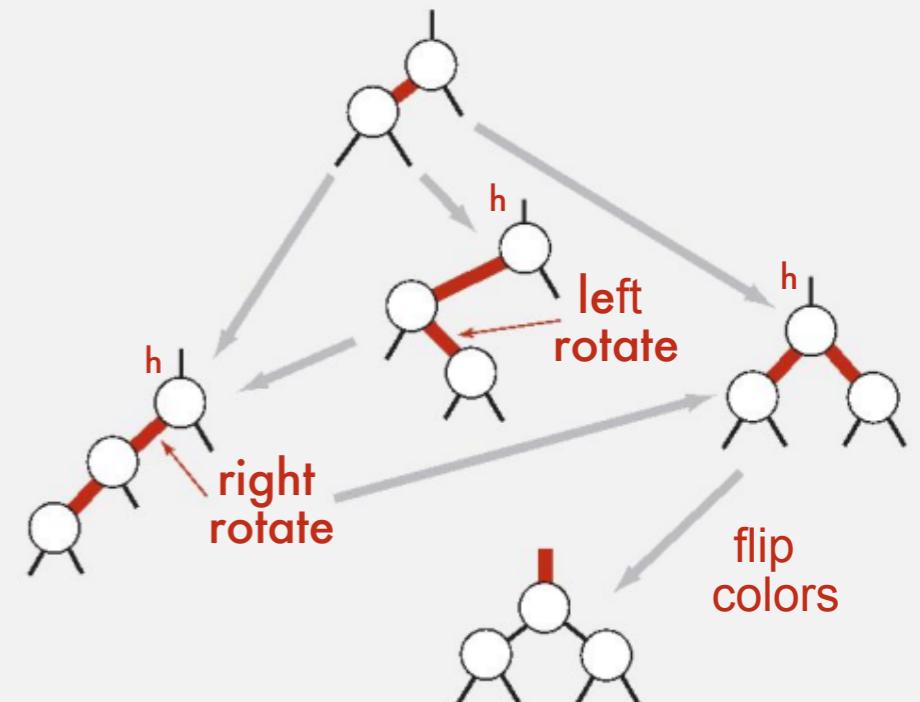
2-3 tree



# Insertion in a LLRB tree: Java implementation

## Same code for all cases.

- Right child red, left child black: **rotate left**.
  - Left child, left-left grandchild red: **rotate right**.
  - Both children red: **flip colors**.



```
private Node put(Node h, Key key, Value val)
{
    if (h == null) return new Node(key, val, RED);
    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 if (cmp == 0) 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)) flipColors(h);

    return h;
}
```

only a few extra lines of code provided

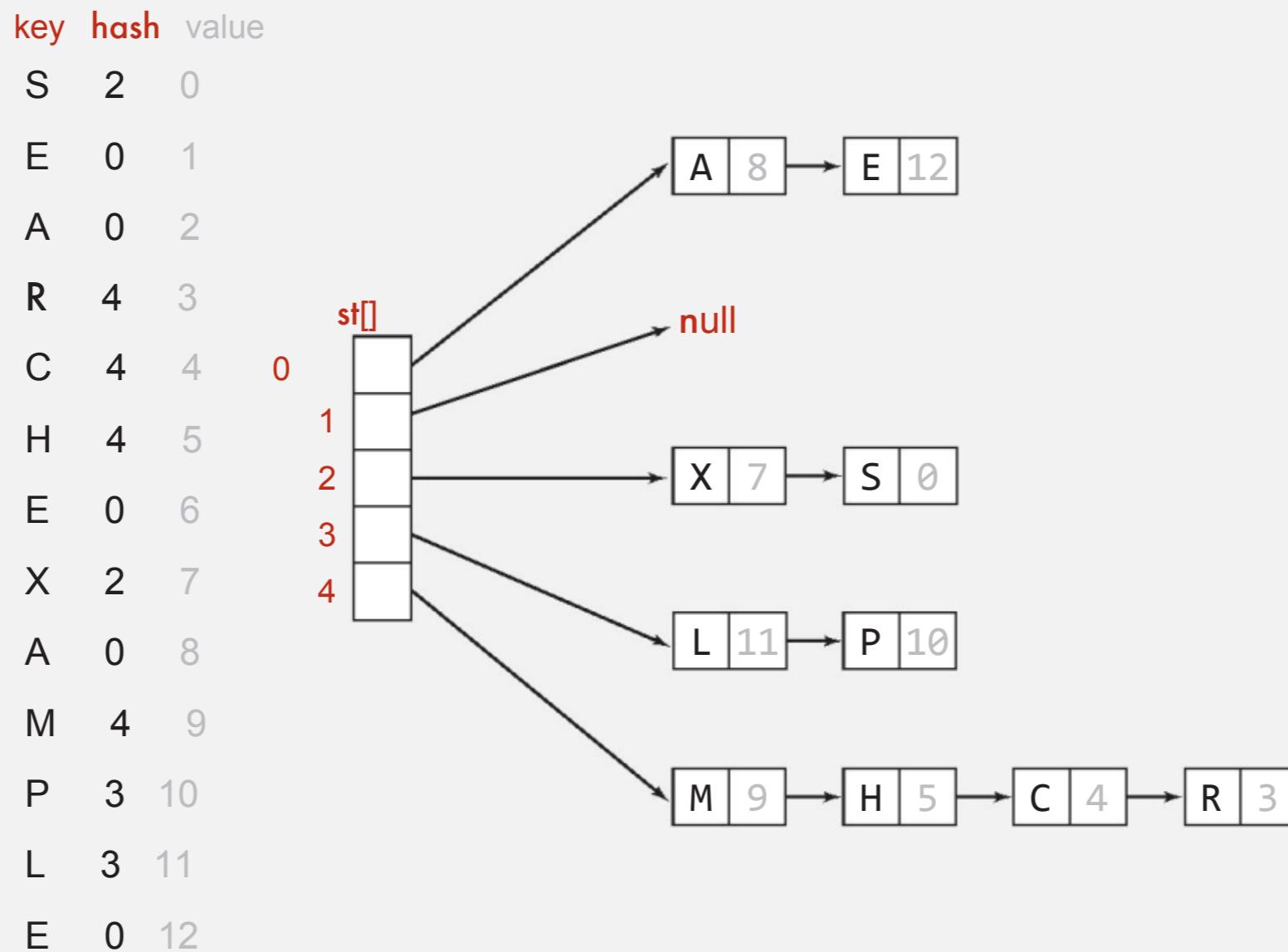
only a few extra lines of code provides near-perfect balance

- insert at bottom  
(and color it red)
- lean left
- balance 4-node
- split 4-node

# Separate-chaining symbol table

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  $i$  chain (if not already there).
- Search: need to search only  $i$  chain.



## Linear-probing hash table summary

---

**Hash.** Map key to integer  $i$  between 0 and  $M-1$ .

**Insert.** Put at table index  $i$  if free; if not try  $i+1, i+2, \dots$

**Search.** Search table index  $i$ ; if occupied but no match, try  $i+1, i+2, \dots$

**Note.** Array size  $M$  must be greater than number of key-value pairs  $N$ .

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
st[]	P	M			A	C	S	H	L		E				R	X

$M = 16$

# ST implementations: summary

implementation	guarantee			average case			ordered ops?	key interface
	search	insert	delete	search hit	insert	delete		
sequential search (unordered list)	N	N	N	$\frac{1}{2}N$	N	$\frac{1}{2}N$		equals()
binary search (ordered array)	$\lg N$	N	N	$\lg N$	$\frac{1}{2}N$	$\frac{1}{2}N$	✓	compareTo()
BST	N	N	N	$1.39 \lg N$	$1.39 \lg N$	$\sqrt{N}$	✓	compareTo()
red-black BST	$2 \lg N$	$2 \lg N$	$2 \lg N$	$1.0 \lg N$	$1.0 \lg N$	$1.0 \lg N$	✓	compareTo()
separate chaining	N	N	N	3-5 *	3-5 *	3-5 *		equals() hashCode()
linear probing	N	N	N	3-5 *	3-5 *	3-5 *		equals() hashCode()

\* under uniform hashing assumption

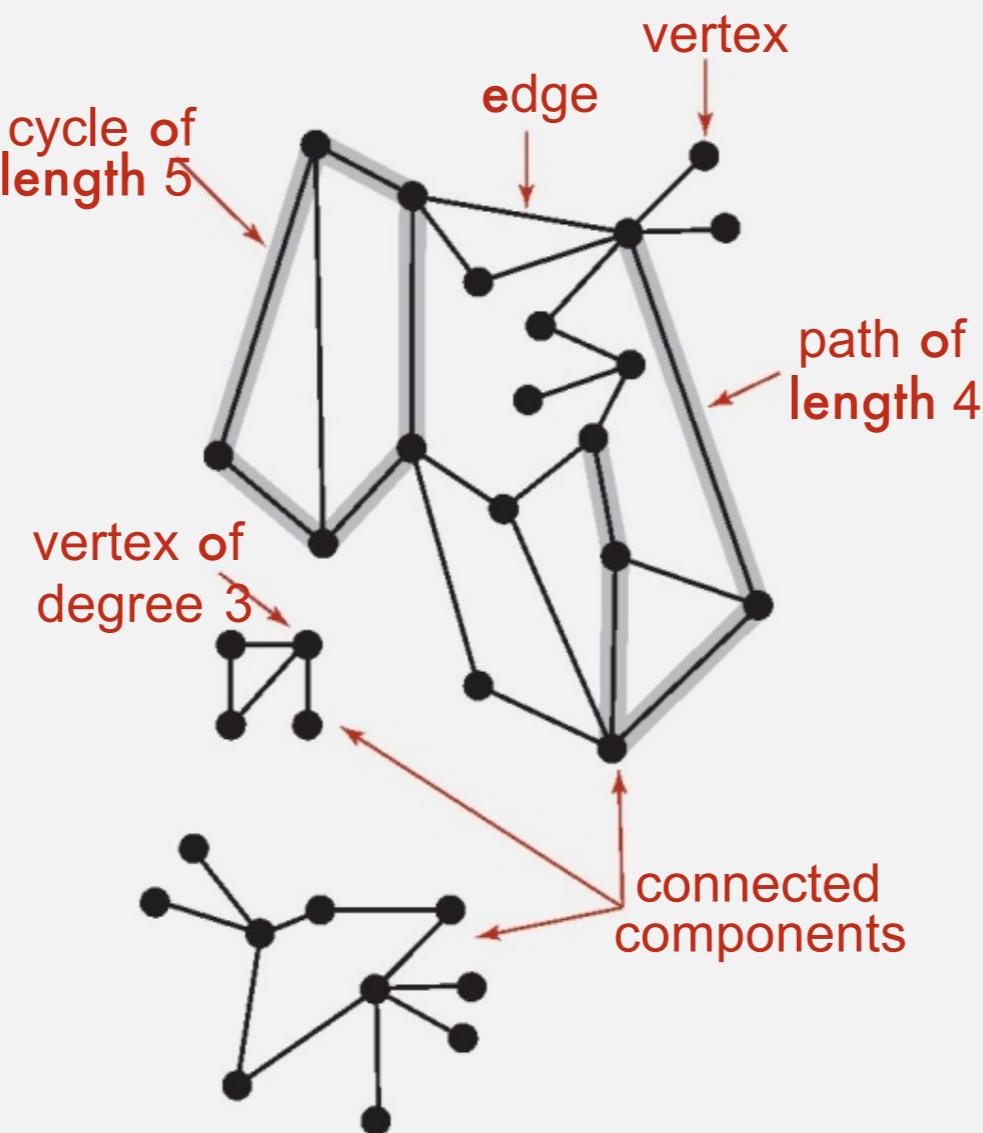
# Graph terminology

---

**Path.** Sequence of vertices connected by edges.

**Cycle.** Path whose first and last vertices **are the same**.

Two vertices **are connected** if there is a path between them.



# Depth-first search

---

Goal. Systematically traverse a graph.

Idea. Mimic maze exploration. ← function-call stack acts as ball of string

DFS (to visit a vertex  $v$ )

Mark  $v$  as visited.

Recursively visit all unmarked  
vertices  $w$  adjacent to  $v$ .

Typical applications.

- Find all vertices connected to a given source vertex.
- Find a path between two vertices.

Design challenge. How to implement?

# Design pattern for graph processing

**Design pattern.** Decouple graph data type from graph processing.

- Create a Graph object.
- Pass the Graph to a graph-processing routine.
- Query the graph-processing routine for information.

```
public class Paths
```

```
    Paths(Graph G, int s)
```

find paths in G from source s

```
    boolean hasPathTo(int v)
```

is there a path from s to v?

```
    Iterable<Integer> pathTo(int v)
```

path from s to v; null if no such path

```
Paths paths = new Paths(G, s);
for (int v = 0; v < G.V(); v++)
    if (paths.hasPathTo(v))
        StdOut.println(v);
```

print all vertices  
connected to s

# Depth-first search: data structures

---

To visit a vertex  $v$  :

- Mark vertex  $v$  as visited.
- Recursively visit all unmarked vertices adjacent to  $v$ .

## Data structures.

- Boolean array `marked[]` to mark visited vertices.
- Integer array `edgeTo[]` to keep track of paths.  
 $(\text{edgeTo}[w] == v)$  means that edge  $v-w$  taken to visit  $w$  for first time
- Function-call stack for recursion.

# Depth-first search: Java implementation

```
public class DepthFirstPaths
{
    private boolean[] marked;
    private int[] edgeTo;
    private int s;

    public DepthFirstPaths(Graph G, int s)
    {
        ...
        dfs(G, s);
    }

    private void dfs(Graph G, int v)
    {
        marked[v] = true;
        for (int w : G.adj(v))
            if (!marked[w])
            {
                dfs(G, w);
                edgeTo[w] = v;
            }
    }
}
```

marked[v]= true  
if v connected to s

edgeTo[v]= previous vertex on path from s to v

initialize data structures

find vertices connected to s

recursive DFS does the work

# Breadth-first search

---

Repeat until queue is empty:

- Remove vertex  $v$  from queue.
- Add to queue all unmarked vertices adjacent to  $v$  and mark them.

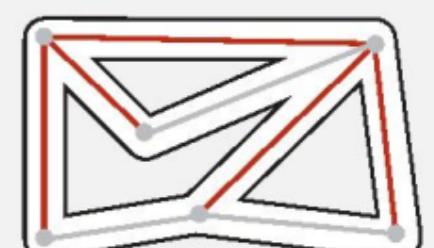
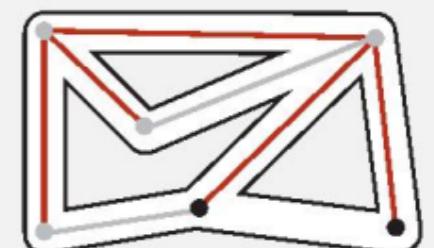
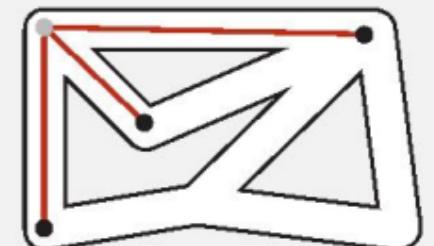
BFS (from source vertex  $s$ )

---

Put  $s$  onto a FIFO queue, and mark  $s$  as visited.

Repeat until the queue is empty:

- remove the least recently added vertex  $v$
  - add each of  $v$ 's unvisited neighbors to the queue,  
and mark them as visited.
- 



# Finding connected components with DFS

```
public class CC
{
    private boolean[] marked;
    private int[] id;
    private int count;
```

$\text{id}[v]$  = id of component containing  $v$   
number of components

```
public CC(Graph G)
{
    marked = new boolean[G.V()];
    id = new int[G.V()];
    for (int v = 0; v < G.V(); v++)
    {
        if (!marked[v])
        {
            dfs(G, v);
            count++;
        }
    }
}
```

run DFS from **one** vertex in  
**each** component

```
public int count()
public int id(int v)
public boolean connected(int v, int w)
private void dfs(Graph G, int v)
```

see next slide

## Finding connected components with DFS (continued)

```
public int count()  
{  return count;  }
```

number of components

```
public int id(int v)  
{  return id[v];  }
```

id of component containing v

```
public boolean connected(int v, int w)  
{ return id[v] == id[w];  }
```

v and w connected iff same id

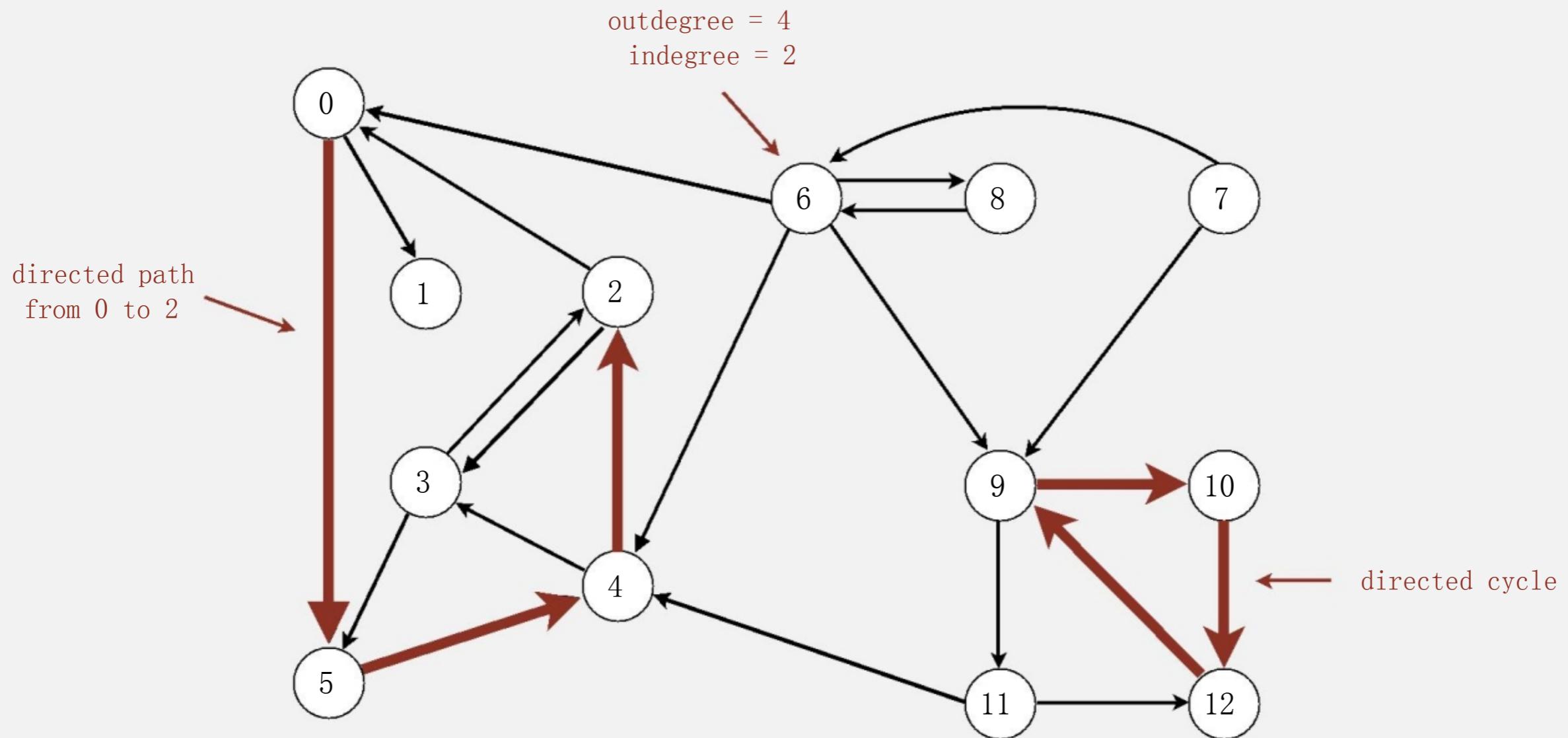
```
private void dfs(Graph G, int v)  
{  
    marked[v] = true;  
    id[v] = count;  
    for (int w : G.adj(v))  
        if (!marked[w])  
            dfs(G, w);  
}
```

all vertices discovered in  
same call of dfs have same id

# Directed graphs

---

Digraph. Set of vertices connected pairwise by directed edges.



# Digraph representations

---

In practice. Use adjacency-lists representation.

- Algorithms based on iterating over vertices pointing from  $v$ .
- Real-world digraphs tend to be sparse.

huge number of vertices,  
small average vertex degree

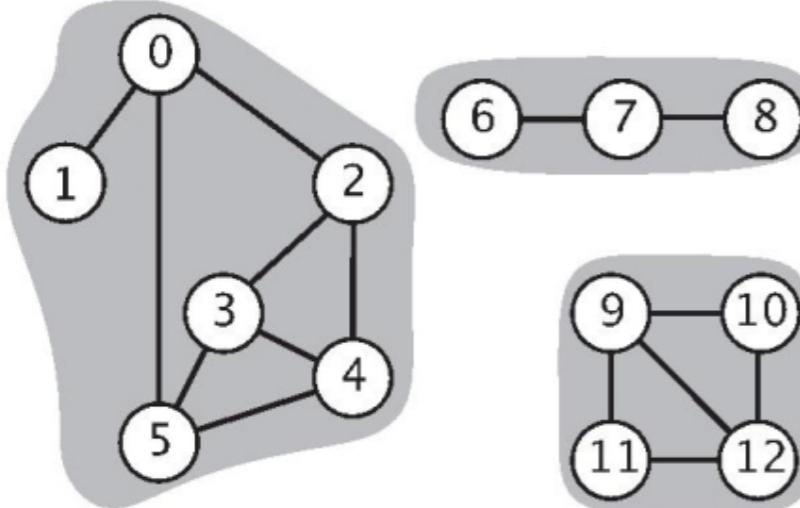
representation	space	insert edge from $v$ to $w$	edge from $v$ to $w$ ?	iterate over vertices pointing from $v$ ?
list of edges	$E$	1	$E$	$E$
adjacency matrix	$V^2$	$1^\dagger$	1	$V$
adjacency lists	$E + V$	1	$\text{outdegree}(v)$	$\text{outdegree}(v)$

<sup>†</sup> disallows parallel edges

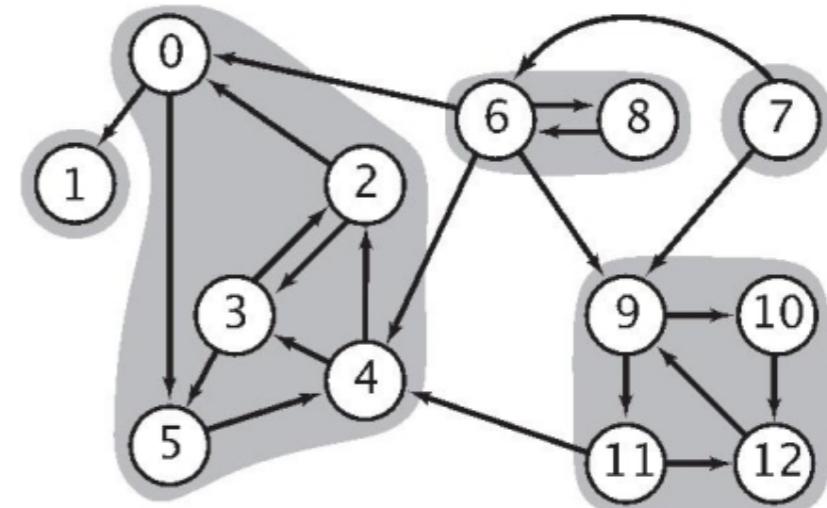
# Connected components vs. strongly-connected components

v and w are **connected** if there is a path between v and w

v and w are **strongly connected** if there is both a directed path from v to w and a directed path from w to v



3 connected components



5 strongly-connected components

connected component id (easy to compute with DFS)

	0	1	2	3	4	5	6	7	8	9	10	11	12
id[]	0	0	0	0	0	0	1	1	1	2	2	2	2

```
public boolean connected(int v, int w)
{ return id[v] == id[w]; }
```

strongly-connected component id (how to compute?)

	0	1	2	3	4	5	6	7	8	9	10	11	12
id[]	1	0	1	1	1	1	3	4	3	2	2	2	2

constant-time client connectivity query

```
public boolean stronglyConnected(int v, int w)
{ return id[v] == id[w]; }
```

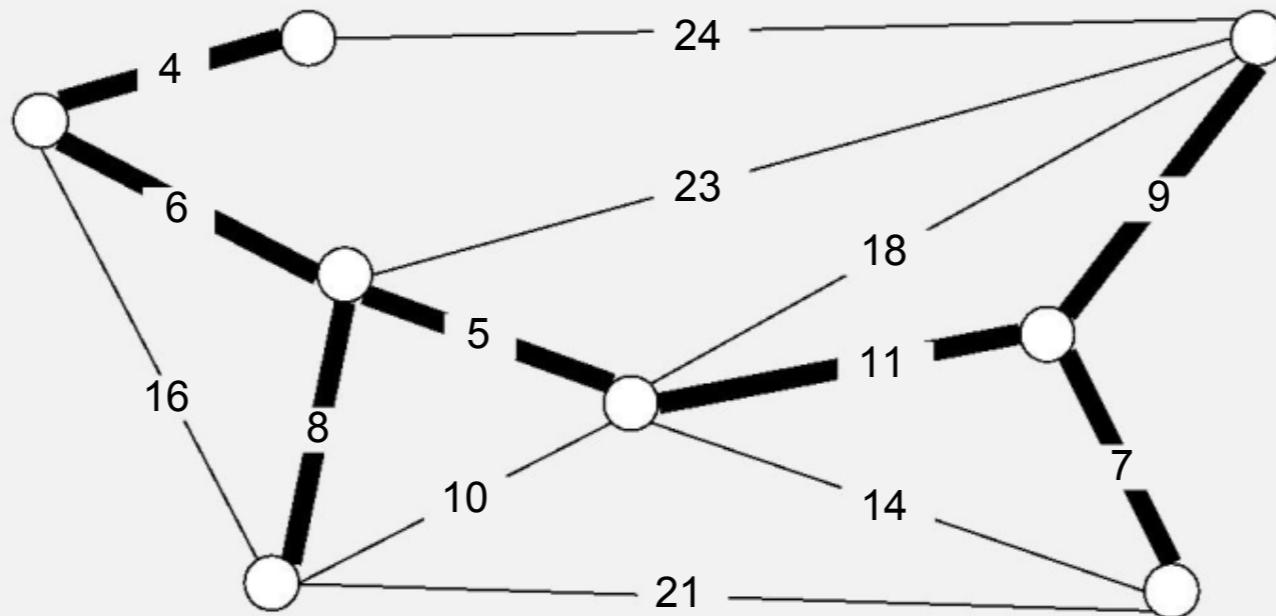
constant-time client strong-connectivity query

# Minimum spanning tree

---

**Given.** Undirected graph  $G$  with positive edge weights (connected).

**Goal.** Find a min weight spanning tree.



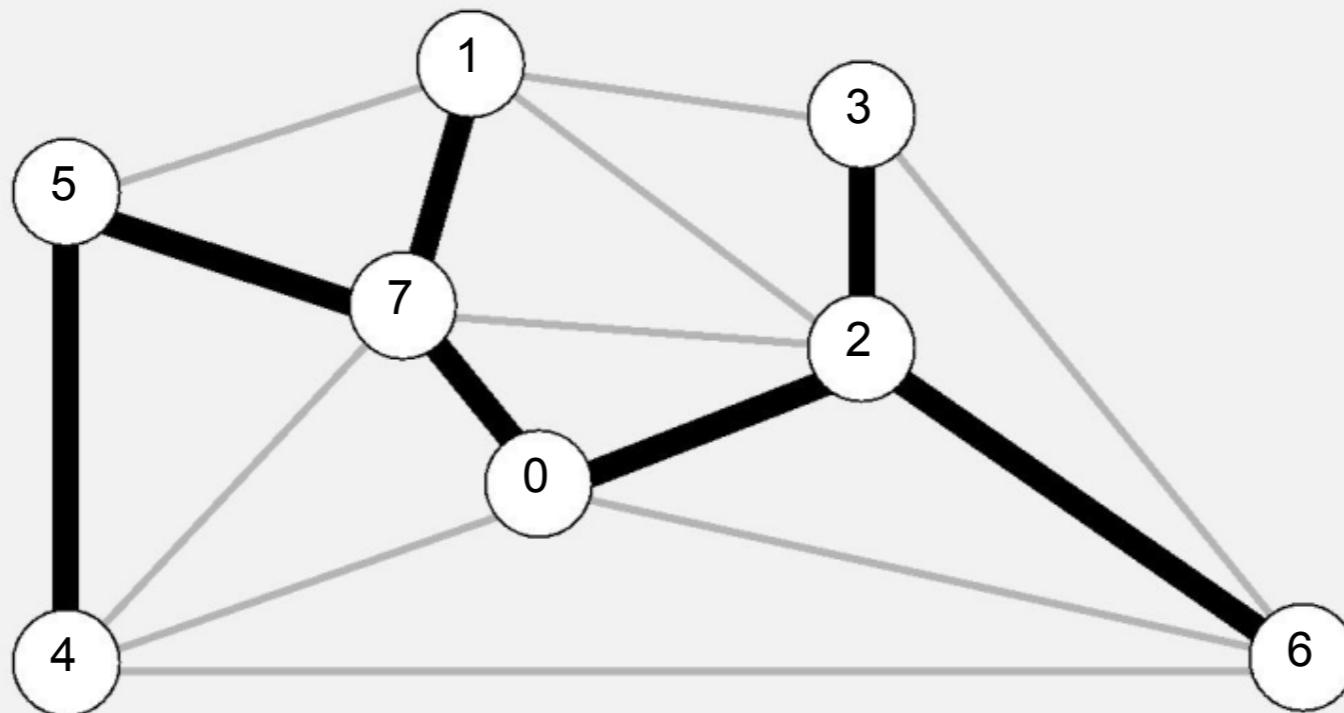
minimum spanning tree  $T$   
(cost =  $50 = 4 + 6 + 8 + 5 + 11 + 9 + 7$ )

**Brute force.** Try all spanning trees?

# Greedy MST algorithm demo

---

- Start with all edges colored gray.
- Find cut with no black crossing edges; color its min-weight edge black.
- Repeat until  $V - 1$  edges are colored black.



MST edges

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

# Key-indexed counting demo

Goal. Sort an array  $a[]$  of  $N$  integers between 0 and  $R - 1$ .

- Count frequencies of each letter using key as index.
- Compute frequency cumulates which specify destinations.
- Access cumulates using key as index to move items.
- Copy back into original array.

```
int N = a.length;
int[] count = new int[R+1];

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

for (int r = 0; r < R; r++)
    count[r+1] += count[r];

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

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

copy  
back

i	a[i]	i	aux[i]
0	a	0	a
1	a	1	a
2	b	2	b
3	b	3	b
4	b	4	b
5	c	5	c
6	d	6	d
7	d	7	d
8	e	8	e
9	f	9	f
10	f	10	f
11	f	11	f

## Key-indexed counting: analysis

---

Proposition. Key-indexed takes time proportional to  $N + R$ .

Proposition. Key-indexed counting uses extra space proportional to  $N + R$ .

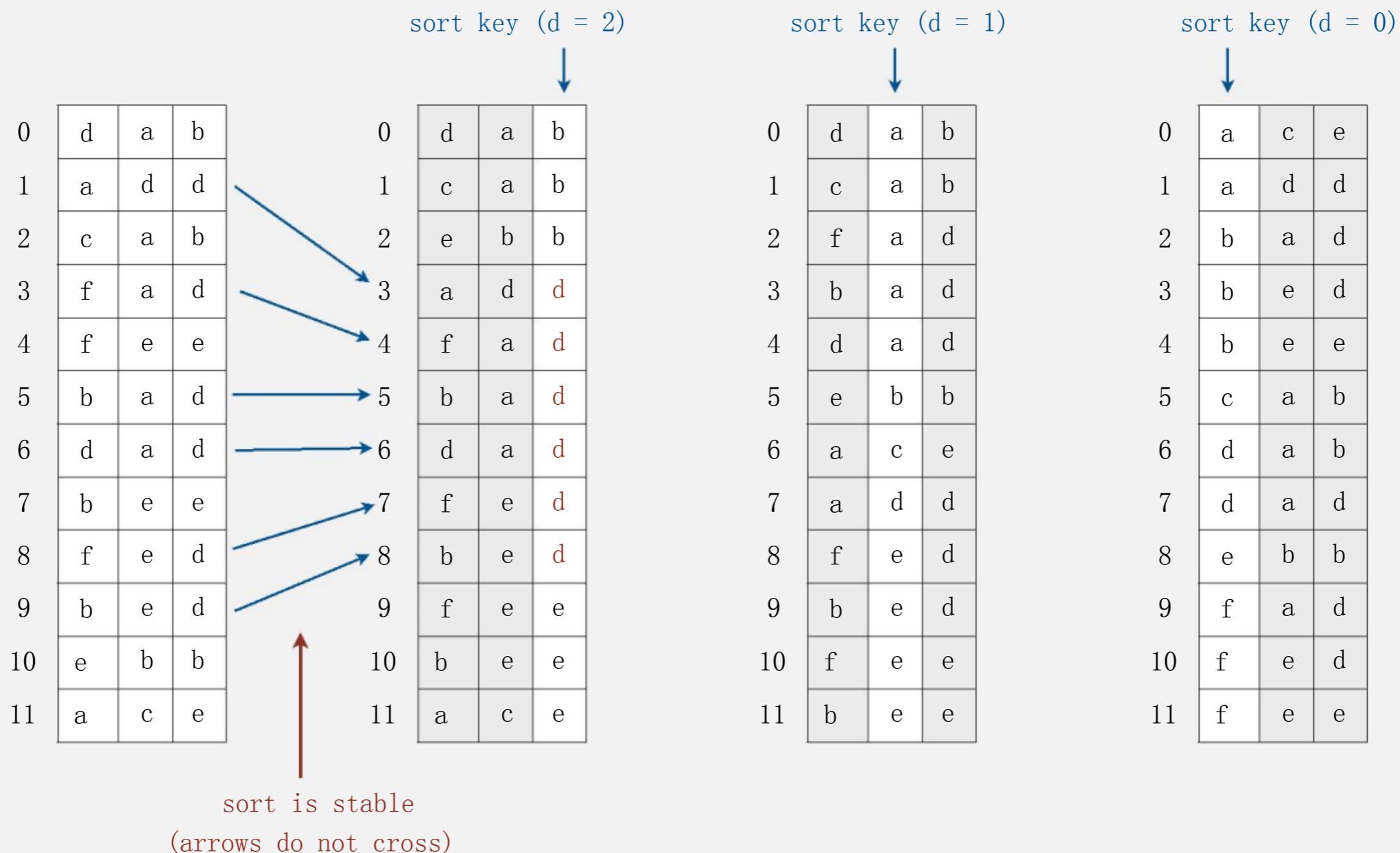
Stable? ✓

a[0]	Anderson	2	Harris	1	aux[0]
a[1]	Brown	3	Martin	1	aux[1]
a[2]	Davis	3	Moore	1	aux[2]
a[3]	Garcia	4	Anderson	2	aux[3]
a[4]	Harris	1	Martinez	2	aux[4]
a[5]	Jackson	3	Miller	2	aux[5]
a[6]	Johnson	4	Robinson	2	aux[6]
a[7]	Jones	3	White	2	aux[7]
a[8]	Martin	1	Brown	3	aux[8]
a[9]	Martinez	2	Davis	3	aux[9]
a[10]	Miller	2	Jackson	3	aux[10]
a[11]	Moore	1	Jones	3	aux[11]
a[12]	Robinson	2	Taylor	3	aux[12]
a[13]	Smith	4	Williams	3	aux[13]
a[14]	Taylor	3	Garcia	4	aux[14]
a[15]	Thomas	4	Johnson	4	aux[15]
a[16]	Thompson	4	Smith	4	aux[16]
a[17]	White	2	Thomas	4	aux[17]
a[18]	Williams	3	Thompson	4	aux[18]
a[19]	Wilson	4	Wilson	4	aux[19]

# Least-significant-digit-first string sort

LSD string (radix) sort.

- Consider characters from right to left.
- Stably sort using  $\text{d}^{\text{th}}$  character as the key (using key-indexed counting).



# Most-significant-digit-first string sort

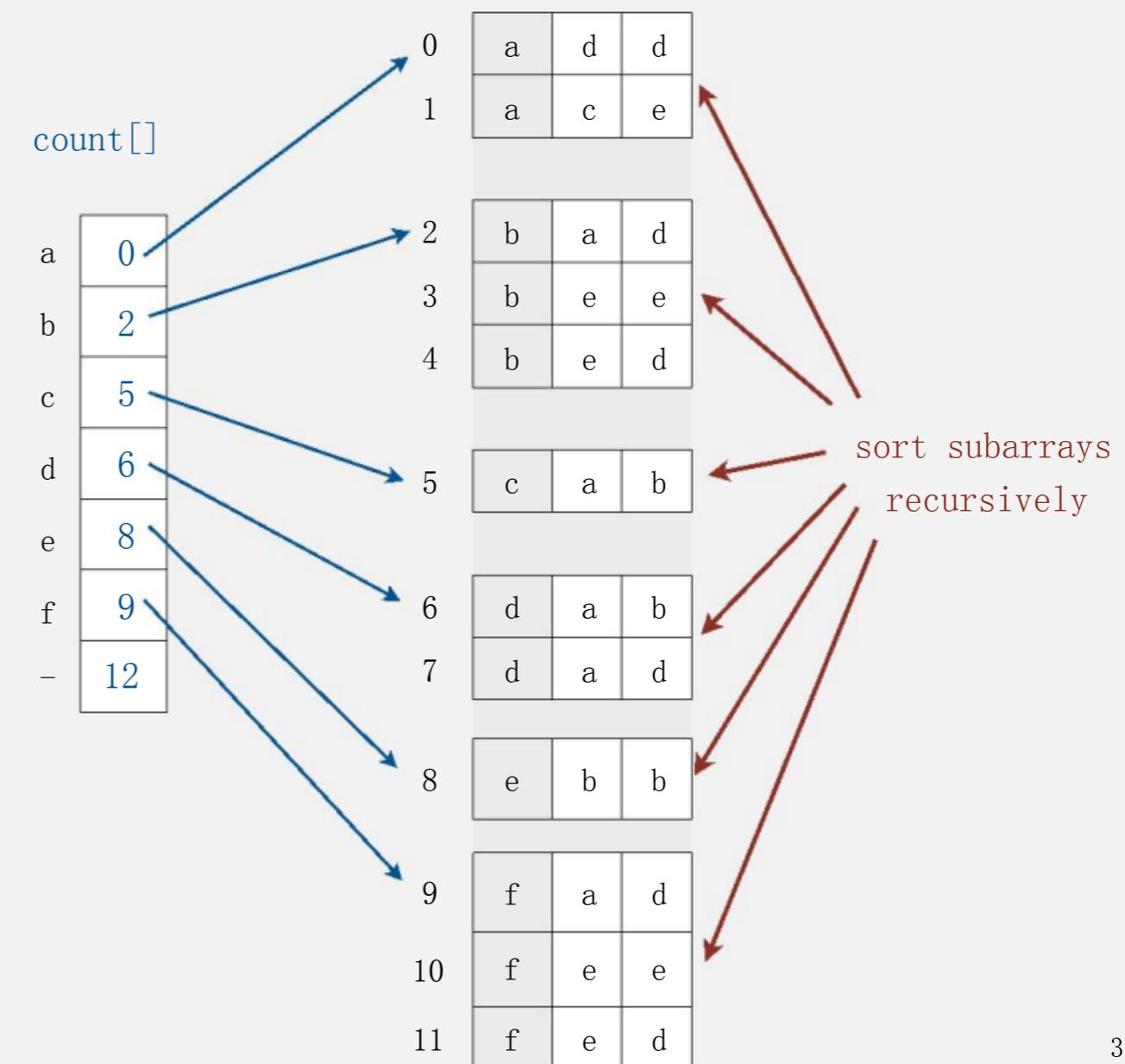
## MSD string (radix) sort.

- Partition array into R pieces according to first character (use key-indexed counting).
- Recursively sort all strings that start with each character (key-indexed counts delineate subarrays to sort).

0	d	a	b
1	a	d	d
2	c	a	b
3	f	a	d
4	f	e	e
5	b	a	d
6	d	a	d
7	b	e	e
8	f	e	d
9	b	e	d
10	e	b	b
11	a	c	e

0	a	d	d
1	a	c	e
2	b	a	d
3	b	e	e
4	b	e	d
5	c	a	b
6	d	a	b
7	d	a	d
8	e	b	b
9	f	a	d
10	f	e	e
11	f	e	d

sort key



## Variable-length strings

Treat strings as if they had an extra char at end (smaller than any char).

0	s	e	a	-1						
1	s	e	a	s	h	e	l	l	s	-1
2	s	e	l	l	s	-1				
3	s	h	e	-1						
4	s	h	e	-1						
5	s	h	e	l	l	s	-1			
6	s	h	o	r	e	-1				
7	s	u	r	e	l	y	-1			

```
private static int charAt(String s, int d)
{
    if (d < s.length()) return s.charAt(d);
    else return -1;
}
```

C strings. Have extra char '\0' at end  $\Rightarrow$  no extra work needed.

# Summary of the performance of sorting algorithms

---

Frequency of operations.

algorithm	guarantee	random	extra space	stable?	operations on keys
insertion sort	$\frac{1}{2}N^2$	$\frac{1}{4}N^2$	1	✓	compareTo()
mergesort	$N \lg N$	$N \lg N$	$N$	✓	compareTo()
quicksort	$1.39 N \lg N$ *	$1.39 N \lg N$	$c \lg N$		compareTo()
heapsort	$2 N \lg N$	$2 N \lg N$	1		compareTo()
LSD sort †	$2 W (N + R)$	$2 W (N + R)$	$N + R$	✓	charAt()
MSD sort ‡	$2 W (N + R)$	$N \log_R N$	$N + D R$	✓	charAt()

D = function-call stack depth  
(length of longest prefix match)



\* probabilistic

† fixed-length W keys

‡ average-length W keys