

Welcome to CS61B!

- Please see the rather extensive information on sections, Covid-19 policy, sections, labs, initial assignments, and the presemester survey on the Spring 2022 CS61B Piazza site.
- Labs start today. In (or preferably before) lab this week, get a CS61B Unix account from <https://inst.eecs.berkeley.edu/webacct>.
- Try logging in remotely to one of the instructional servers (...@ X .berkeley.edu, where X is ashby.cs, derby.cs, cedar.cs, cory.eecs, and others).
- The course homepage (<https://inst.eecs.berkeley.edu/cs61b/sp22>) is our central distribution site for assignments, lecture slides, course policy, and much else.
- Lectures will be recorded and screencast. The recordings should become available in the bCourses Media Gallery sometime after the lecture.

Crowding, etc.

- If you choose not to take this course please drop it as soon as possible for the benefit of others (the add/drop deadline is 9 February—28 January if you wish to avoid a fee).
- As you know, Dwinelle will not hold us all, which is why there are both offline and online lectures. Lecture seating is on a first-come-first-seated basis. Definitely not ideal, but we hope that after the first few weeks, those of you who prefer in-person lectures will be able to have it.
- Lectures, etc., will be entirely online for at least the first two weeks, due to the current outbreak.

Texts

- There are two readers currently on-line (see the website).
- Textbook (for first part of the course only) is *Head First Java*. It's kind of silly, but has the necessary material.

Course Organization I

- You read; we illustrate.
- Labs are important: exercise of programming principles as well as practical dirty details go there. Generally we will give you homework points for doing them.
- Homework is important, but it's reasonably easy to get full credit: use it as you see fit and *turn it in!* You should get points for just putting some reasonable effort into it.
- Individual projects are *really* important! Expect to learn a lot. Projects are *not* team efforts (that's for later courses).

Course Organization II

- Use of tools *is* part of the course. Programming takes place in a *programming environment*:
 - Handles editing, debugging, compilation, archiving versions.
 - Personally, I keep it simple: Emacs + gjdb + make + git, (documented in one of the readers and on-line). But we'll look at IntelliJ in lab.
- Tests are challenging: better to stay on top than to cram.
- Tests, 40%; Projects, 50%; HW, 10%
- Stressed? Tell us!

Pandemic Considerations

- It's everyone's responsibility to look out for each other.
- This semester, in particular, this means adhering to certain inconvenient practices mandated by the University.
- These include wearing masks indoors, as well as staying home when sick.
- Please observe the mask mandate; if anyone refuses, I can and will be forced to simply end the day's lecture, and you'll all have to rely on the on-line slides for the material.

Academic Dishonesty

- Sadly, the incidence of academic dishonesty seems to have increased over the years.
- To an extent, this is our fault: the minimum GPA threshold policy for L&S majors puts people under a lot of stress,
- Nevertheless, we can't afford to tolerate cheating. The Course Info tab on the course homepage contains our policy on cheating and the penalties we impose; please read them.
- By keeping up with the course and starting assignments early, you can reduce any perceived need to cheat.
- Also, this course is not curved, so you are not disadvantaged by other people's dishonesty.

Programming, not Java

- Here, we learn *programming*, not Java (or Unix, or Windows, or...)
- Programming principles span many languages
 - Look for connections.
 - Syntax ($x+y$ vs. $(+ x y)$) is superficial.
 - Java, Python, and Scheme have a lot in common.
- Whether you use GUIs, text interfaces, or embedded systems, important ideas are the same.

For next time

- Please read Chapter 1 of *Head First Java*, plus §1.1-1.9 of the on-line book *A Java Reference*, available on the class website.
- This is an overview of most of Java's features.
- We'll start looking at examples on Friday.
- Always remember the questions that come up when you read something we assign:
 - Who knows? We might have made a mistake.
 - Feel free to ask at the start of lectures, by email, or by Piazza.

Acronyms of Wisdom

DBC

RTFM

A Quick Tour through the First Program

In Python, we would write

```
# Traditional first program
print("Hello, world")
```

But in Java,

```
/** Traditional first program.
 * @author P. N. Hilfinger */
public class Hello {
    /** Print greeting. ARGS is ignored. */
    public static void main(String[] args) {
        System.out.println("Hello, world!");
    }
}
```

Commentary

```
/** Traditional first program.  
 * @author P. N. Hilfinger */  
public class Hello {  
    /** Print greeting.  ARGS is ignored.  */  
    public static void main(String[] args) {  
        System.out.println("Hello, world!");  
    }  
}
```

- Java comments can either start with '//' and go to the end of the line (like '#' in Python), or they can extend over any number of lines, bracketed by '/*' and '*/'.
- I don't use the '//' comments, except for things that are supposed to be replaced, and our style checks will flag them.
- The second, multiline kind of comment includes those that start with '/**', which are called *documentation comments* or *doc comments*.
- Documentation comments are just comments, having no effect, but various tools interpret them as providing documentation for the things that follow them. They're generally a good idea and our style checks require them.

Classes

```
/** Traditional first program.  
 * @author P. N. Hilfinger */  
public class Hello {  
    /** Print greeting. ARGS is ignored. */  
    public static void main(String[] args) {  
        System.out.println("Hello, world!");  
    }  
}
```

- Every function and variable in Java is contained in some *class*.
- These are like Python's classes, but with (of course) numerous differences in detail.
- All classes, in turn, belong to some *package*. The Hello class belongs to the *anonymous package*.
- We'll see named packages later,

Methods (Functions)

```
/** Traditional first program.  
 * @author P. N. Hilfinger */  
public class Hello {  
    /** Print greeting. ARGS is ignored. */  
    public static void main(String[] args) {  
        System.out.println("Hello, world!");  
    }  
}
```

- Function headers in Java contain more information than those in Python. They specify the *types* of values *returned* by the function and taken as *parameters* to the functions.
- The “type” void has no possible values; the *main* function here returns nothing. The type String is like Python’s str. The trailing ‘[]’ means *array of*. Arrays are like Python lists, except that their size is fixed once created.
- Hence, *main* takes a list of strings and returns nothing.
- Functions named “main” and defined like the example above are special: they are what get called when one runs a Java program (in Python, the main function is essentially anonymous).

Selection

```
/** Traditional first program.  
 *  @author P. N. Hilfinger */  
public class Hello {  
    /** Print greeting. ARGS is ignored. */  
    public static void main(String[] args) {  
        System.out.println("Hello, world!");  
    }  
}
```

- As in Python, $\mathcal{E}.N$ means “the thing named N that is in or that applies to the thing identified (or computed) by \mathcal{E} .”
- Thus “System.out” means “the variable named ‘out’ that is found in the class named ‘System’.”
- Likewise, “System.out.println” means “the method named ‘println’ that applies to the object referenced by the value of variable ‘System.out’.”

Access

```
/** Traditional first program.  
 * @author P. N. Hilfinger */  
public class Hello {  
    /** Print greeting. ARGS is ignored. */  
    public static void main(String[] args) {  
        System.out.println("Hello, world!");  
    }  
}
```

- Every declared entity in Java has *access permissions* indicating what pieces of code may mention it.
- In particular, *public* classes, methods, and variables may be referred to anywhere else in the program.
- We sometimes refer to them as *exported* from their class (for methods or variables) or package (for classes).

Static

```
/** Traditional first program.  
 * @author P. N. Hilfinger */  
public class Hello {  
    /** Print greeting. ARGS is ignored. */  
    public static void main(String[] args) {  
        System.out.println("Hello, world!");  
    }  
}
```

- Static methods and variables are “one-of” things.
- A static method is just like an ordinary Python function (outside of any class) or a function in a Python class that is annotated `@staticmethod`.
- A static variable is like a Python variable defined outside of any class or a variable selected from a class, as opposed to from a class instance.
- Other variables are local variables (in functions) or instance variables (in classes), and these are as in Python.

Administrivia

- Please make sure you have obtained a Unix account.
- If you decide not to take this course after all, please tell CalCentral ASAP, so that we can adjust the waiting list accordingly.
- HW #0 will be due next Friday at midnight. While you get credit for any submission, we *strongly* suggest that you give the problems a serious try.
- We *strongly discourage* taking this course P/NP (or S/U).

Lecture #2: Let's Write a Program: Prime Numbers

Problem: want java Primes U to print prime numbers through U .

You type: java Primes 101

It types: 2 3 5 7 11 13 17 19 23 29
31 37 41 43 47 53 59 61 67 71
73 79 83 89 97 101

Definition: A **prime** number is an integer greater than 1 that has no divisors smaller than itself other than 1.

(Alternatively: $p > 1$ is prime iff $\gcd(p, x) = 1$ for all $0 < x < p$.)

Plan

```
public class Primes {  
    /** Print all primes up to ARGS[0] (interpreted as an  
     * integer), 10 to a line. */  
    public static void main(String[] args) {  
        printPrimes(Integer.parseInt(args[0]));  
    }  
  
    /** Print all primes up to and including LIMIT, 10 to  
     * a line. */  
    private static void printPrimes(int limit) {  
        /*{ For every integer, x, between 2 and LIMIT, print it if  
           isPrime(x), 10 to a line. }*/  
    }  
  
    /** True iff X is prime. */  
    private static boolean isPrime(int x) {  
        return /*( X is prime );*/  
    }  
}
```

Testing for Primes

```
private static boolean isPrime(int x) {  
    return /*( X is prime */;  
}
```

Testing for Primes

```
private static boolean isPrime(int x) {  
    if (x <= 1)  
        return false;  
    else  
        return !isDivisible(x, 2, x); // "!" means "not"  
}
```

Testing for Primes

```
private static boolean isPrime(int x) {  
    if (x <= 1)  
        return false;  
    else  
        return !isDivisible(x, 2, x); // "!" means "not"  
}  
  
/** True iff X is divisible by any positive number >= LOW >= 1  
 * and < HIGH. */  
private static boolean isDivisible(int x, int low, int high) {  
    return /*( True iff x is divisible by k, low<=k<high. );*/;  
}
```

Testing for Primes

```
private static boolean isPrime(int x) {  
    if (x <= 1)  
        return false;  
    else  
        return !isDivisible(x, 2, x); // "!" means "not"  
}  
  
/** True iff X is divisible by any positive number >= LOW >= 1  
 * and < HIGH. */  
private static boolean isDivisible(int x, int low, int high) {  
    if (low >= high) // a "guard"  
        return false;  
    else if (x % low == 0) // "%" means "remainder"  
        return true;  
    else // if (low < high && x % low != 0)  
        return isDivisible(x, low, high);  
}
```

Thinking Recursively

Understand and check `isDivisible(13, 2)` by *tracing one level*.

```
/** True iff X is divisible by some number
 *  >= LOW >= 1 and < HIGH. */
private static boolean isDivisible...
    if (low >= high)
        return false;
    else if (x % low == 0)
        return true;
    else
        return isDivisible(x, low + 1, high);
}
```

Lesson: Comments aid understanding.
Make them *count!*

- Call assigns $x=13$, $\text{low}=2$, $\text{high}=13$
- Body has form
 - if ($\text{low} \geq \text{high}$) S_1 else S_2 .
- Since $2 < 13$, we evaluate the (first) else.
- Check if $13 \bmod 2 = 0$; it's not.
- Left with `isDivisible(13, 3, 13)`.
- Rather than tracing it, instead *use the comment*:
 - Since 13 is *not* divisible by any integer in the range 3..12, `isDivisible(13, 3, 13)` must be *false*, and we're done!
 - Sounds like that last step begs the question. Why doesn't it?

Iteration

- `isDivisible` is *tail recursive*, and so creates an *iterative process*.
- Traditional “Algol family” production languages have special syntax for iteration. Four equivalent versions of `isDivisible`:

```
if (low >= high)
    return false;
else if (x % low == 0)
    return true;
else
    return isDivisible(x, low+1, high); }
```

```
while (low < high) { // !(low >= high)
    if (x % low == 0)
        return true;
    low = low+1;
    // or low += 1, or (yuch) low++
}
return false;
```

```
int k = low;
while (k < high) {
    if (x % k == 0)
        return true;
    k += 1;
}
return false;
```

```
for (int k = low; k < high; k += 1) {
    if (x % k == 0)
        return true;
}
return false;
```

Using Facts about Primes

- A couple of obvious facts:
 - $k \leq \sqrt{N}$ iff $N/k \geq \sqrt{N}$, for $N, k > 0$.
 - If k divides N then N/k divides N .
- So how far do we really have to go to find a possible divisor for x ?

Using Facts about Primes

- A couple of obvious facts:
 - $k \leq \sqrt{N}$ iff $N/k \geq \sqrt{N}$, for $N, k > 0$.
 - If k divides N then N/k divides N .
- So how far do we really have to go to find a possible divisor for x ?
Only up to and including \sqrt{x} .

Using Facts about Primes

- A couple of obvious facts:
 - $k \leq \sqrt{N}$ iff $N/k \geq \sqrt{N}$, for $N, k > 0$.
 - If k divides N then N/k divides N .
- So how far do we really have to go to find a possible divisor for x ?
Only up to and including \sqrt{x} .
- So, reimplement `isPrime`:

```
private static boolean isPrime(int x) {  
    if (x <= 1)  
        return false;  
    else  
        return !isDivisible(x, 2, (int) Math.round(Math.sqrt(x)));  
        // (int) ... here converts to an integer in the range  
        // -231..231 - 1 (type 'int') from one in the  
        // range -263..263 - 1 (type 'long').  
}
```

Cautionary Aside: Floating Point

- In the last slide, we used

```
(int) Math.round(Math.sqrt(x));
```

intending that this would check all values of k up to and including the square root of x .

- Since floating-point operations yield *approximations* to the corresponding mathematical operations, you might ask the following about `Math.round(Math.sqrt(x))`:
 - Is it always at least $\lfloor \sqrt{x} \rfloor$? ($\lfloor z \rfloor$ means "the largest integer $\leq z$."
If not, we might miss testing \sqrt{x} when x is a perfect square.
- As it happens, the answer is "yes" for IEEE floating-point square roots.
- Just an example of the sort of detail that must be checked in edge cases.

Final Task: printPrimes (Simplified)

```
/** Print all primes up to and including LIMIT. */  
private static void printPrimes(int limit) {  
  
}
```

Simplified printPrimes Solution

```
/** Print all primes up to and including LIMIT. */
private static void printPrimes(int limit) {
    for (int p = 2; p <= limit; p += 1) {
        if (isPrime(p)) {
            System.out.print(p + " ");
        }
    }
    System.out.println();
}
```

printPrimes (full version)

```
/** Print all primes up to and including LIMIT, 10 to
 * a line. */
private static void printPrimes(int limit) {
    int np;
    np = 0;
    for (int p = 2; p <= limit; p += 1) {
        if (isPrime(p)) {
            System.out.print(p + " ");
            np += 1;
            if (np % 10 == 0)
                System.out.println();
        }
    }
    if (np % 10 != 0)
        System.out.println();
}
```

Recreation

Prove that $\lfloor(2 + \sqrt{3})^n\rfloor$ is odd for all integer $n \geq 0$.

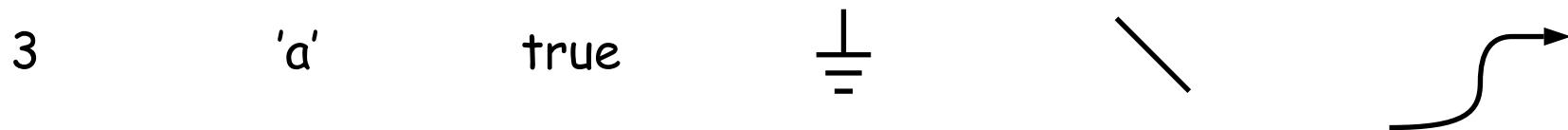
[Source: D. O. Shklarsky, N. N. Chentzov, I. M. Yaglom, *The USSR Olympiad Problem Book*, Dover ed. (1993), from the W. H. Freeman edition, 1962.]

CS61B Lecture #3: Values and Containers

- Labs are normally due at midnight Friday. Last week's lab, however, is due this coming Friday at midnight.
- Today. Simple classes. Scheme-like lists. Destructive vs. non-destructive operations. Models of memory.

Values and Containers

- **Values** are numbers, booleans, and pointers. **Values never change.** (So, for example, the assignment $3 = 2$ would be invalid.)



- **Simple containers** contain values:

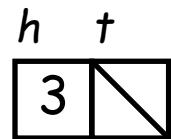


Examples: variables, fields, individual array elements, parameters.
The **contents** of containers can change.

Structured Containers

Structured containers contain (0 or more) other containers:

Class Object



Array Object

0	1	2
42	17	9

Empty Object



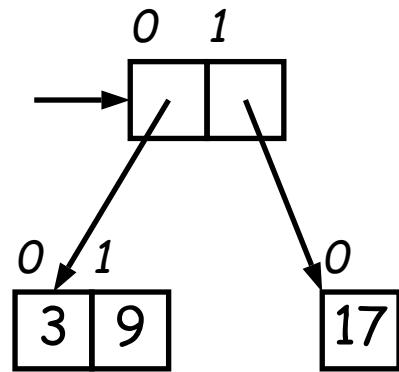
*Alternative
Notation*

<i>h:</i>	3
<i>t:</i>	

0	42
1	17
2	9

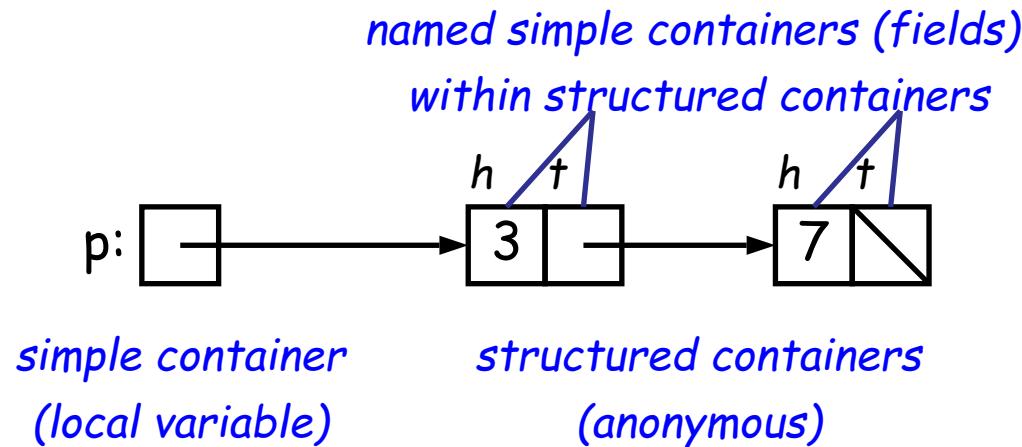
Pointers

- *Pointers* (or *references*) are values that *reference* (point to) containers.
- One particular pointer, called **null**, points to nothing.
- In Java, structured containers contain only simple containers, but pointers allow us to build arbitrarily big or complex structures anyway.



Containers in Java

- Containers may be *named* or *anonymous*.
- In Java, *all* simple containers are named, *all* structured containers are anonymous, and pointers point only to structured containers. (Therefore, structured containers contain only simple containers).



- In Java, assignment copies values into simple containers.
- *Exactly* like Scheme and Python!
- (Python also has slice assignment, as in `x[3:7] = ...`, which is shorthand for something else entirely.)

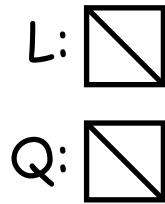
Defining New Types of Object

- Class declarations introduce new types of objects.
- Example: list of integers:

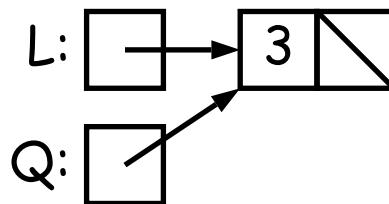
```
public class IntList {  
    // Constructor function (used to initialize new object)  
    /** List cell containing (HEAD, TAIL). */  
    public IntList(int head, IntList tail) {  
        this.head = head; this.tail = tail;  
    }  
  
    // Names of simple containers (fields)  
    // WARNING: public instance variables usually bad style!  
    public int head;  
    public IntList tail;  
}
```

Primitive Operations

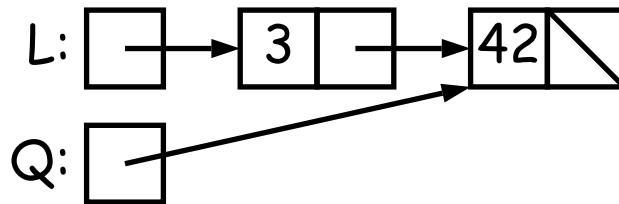
```
IntList Q, L;
```



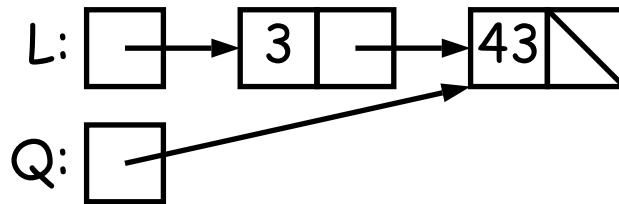
```
L = new IntList(3, null);  
Q = L;
```



```
Q = new IntList(42, null);  
L.tail = Q;
```

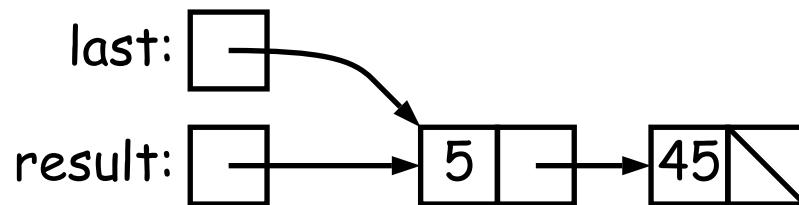


```
L.tail.head += 1;  
// Now Q.head == 43  
// and L.tail.head == 43
```

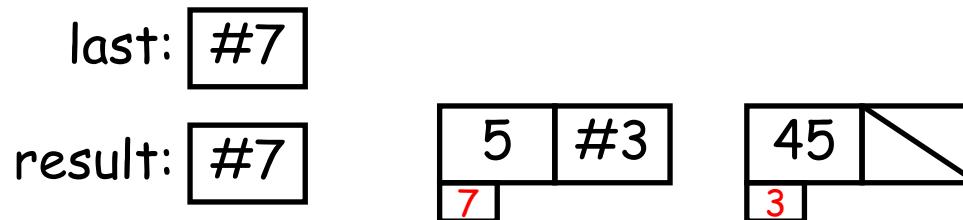


Side Excursion: Another Way to View Pointers

- Some folks find the idea of “copying an arrow” somewhat odd.
- Alternative view: think of a pointer as a *label*, like a street address.
- Each object has a permanent label on it, like the address plaque on a house.
- Then a variable containing a pointer is like a scrap of paper with a street address written on it.
- One view:

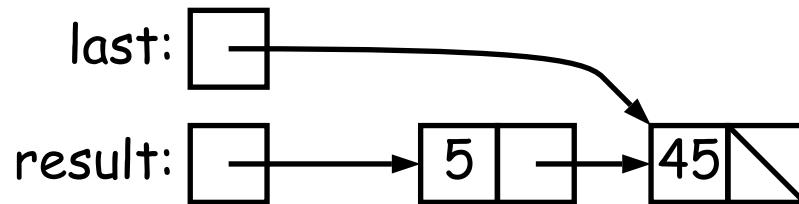


- Alternative view:

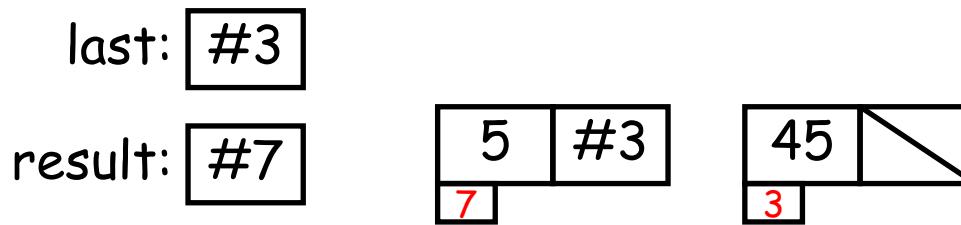


Another Way to View Pointers (II)

- Assigning a pointer to a variable looks just like assigning an integer to a variable.
- So, after executing “`last = last.tail;`” we have



- Alternative view:



- Under alternative view, you might be less inclined to think that assignment would change object #7 itself, rather than just “`last`”.
- BEWARE! Internally, pointers really are just numbers, but Java treats them as more than that: they have *types*, and you can't just change integers into pointers.

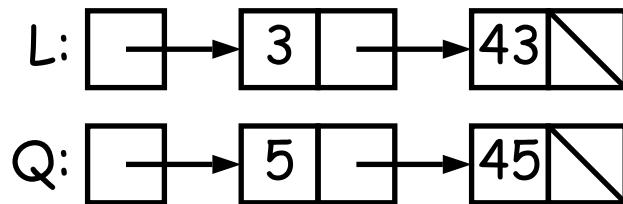
Destructive vs. Non-destructive

Problem: Given a (pointer to a) list of integers, L , and an integer increment n , return a list created by incrementing all elements of the list by n .

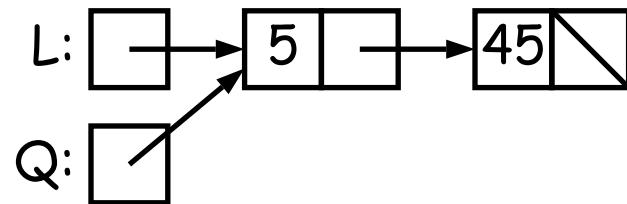
```
/** List of all items in P incremented by n. Does not modify
 * existing IntLists. */
static IntList incrList(IntList P, int n) {
    return /*( P, with each element incremented by n )*/
}
```

We say `incrList` is *non-destructive*, because it leaves the input objects unchanged, as shown on the left. A *destructive* method may modify the input objects, so that the original data is no longer available, as shown on the right:

After $Q = \text{incrList}(L, 2)$:



After $Q = \text{dincrList}(L, 2)$ (destructive):



Nondestructive IncrList: Recursive

```
/** List of all items in P incremented by n. */
static IntList incrList(IntList P, int n) {
    if (P == null) {
        return null;
    } else {
        return new IntList(P.head+n, incrList(P.tail, n));
    }
}
```

- In the call `incrList(P, 2)`, where `P` contains 3 and 43, which `IntList` object gets created first?

Nondestructive IncrList: Recursive

```
/** List of all items in P incremented by n. */
static IntList incrList(IntList P, int n) {
    if (P == null) {
        return null;
    } else {
        return new IntList(P.head+n, incrList(P.tail, n));
    }
}
```

- In the call `incrList(P, 2)`, where `P` contains 3 and 43, which `IntList` object gets created first?
- Answer: The last one.

Nondestructive IncrList: Why Return the Value?

- If I want to update Q to an incremented list, why must I write

```
Q = incrList(Q, 4); ?
```

- Couldn't I instead just write

```
incrList2(Q, 4);
```

and define

```
/** List of all items in P incremented by n. */
static IntList incrList2(IntList P, int n) {
    if (P == null) {
        P = null;
    } else {
        P = new IntList(P.head+n, incrList2(P.tail, n));
    }
    return P;
} ?
```

Nondestructive IncrList: Why Return the Value?

- If I want to update Q to an incremented list, why must I write

```
Q = incrList(Q, 4); ?
```

- Couldn't I instead just write

```
incrList2(Q, 4);
```

and define

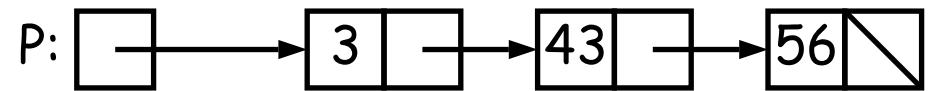
```
/** List of all items in P incremented by n. */
static IntList incrList2(IntList P, int n) {
    if (P == null) {
        P = null;
    } else {
        P = new IntList(P.head+n, incrList2(P.tail, n));
    }
    return P;
} ?
```

- No. Assigning to the formal parameter does not affect the actual. Java uses call by value, just like Python.

An Iterative Version

An iterative incrList is tricky, because it is *not* tail recursive.
Easier to build things first-to-last, unlike recursive version:

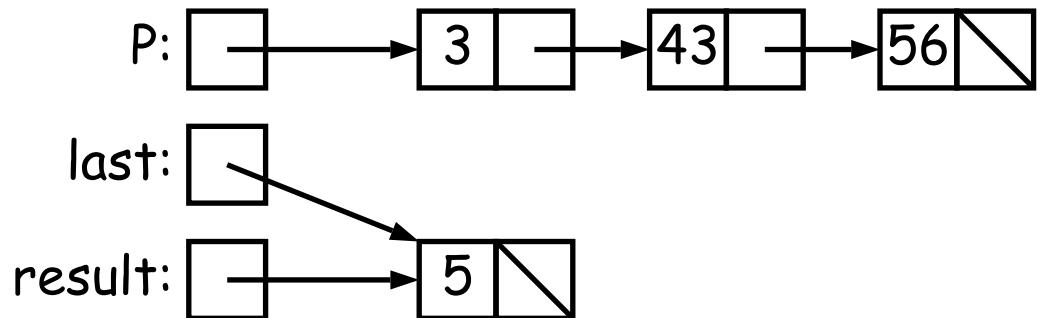
```
static IntList incrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    IntList result, last;  
  
    return result;  
}
```



An Iterative Version

An iterative incrList is tricky, because it is *not* tail recursive.
Easier to build things first-to-last, unlike recursive version:

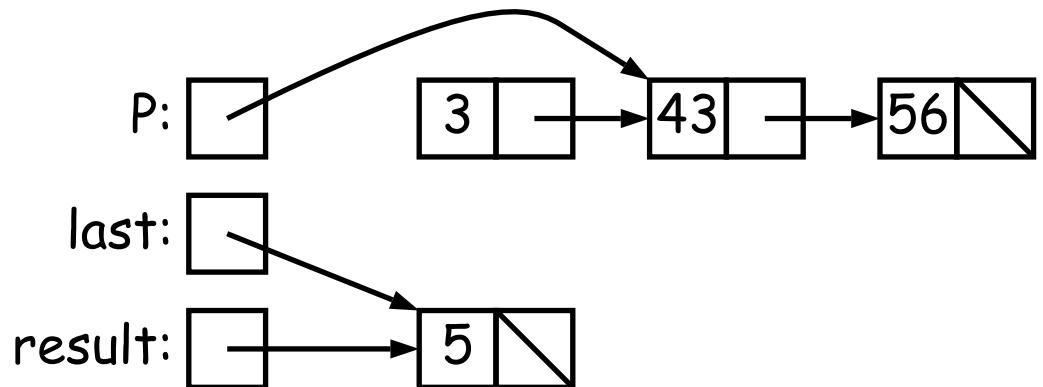
```
static IntList incrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    IntList result, last;  
  
    return result;  
}
```



An Iterative Version

An iterative incrList is tricky, because it is *not* tail recursive.
Easier to build things first-to-last, unlike recursive version:

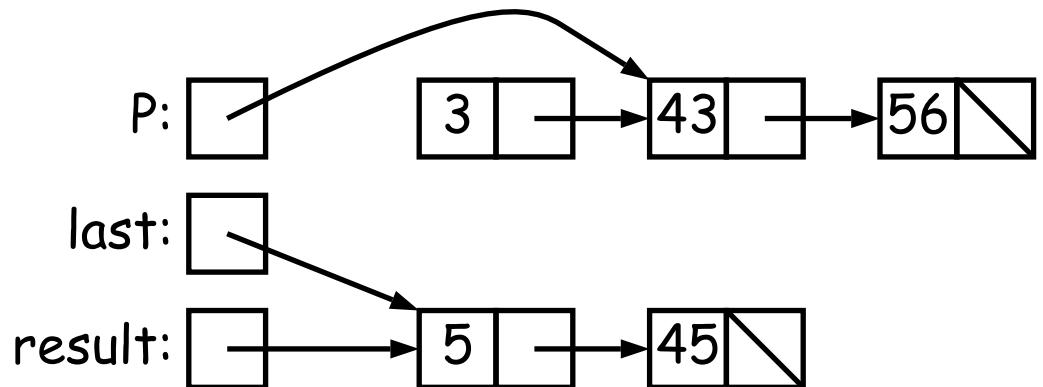
```
static IntList incrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    IntList result, last;  
  
    return result;  
}
```



An Iterative Version

An iterative incrList is tricky, because it is *not* tail recursive.
Easier to build things first-to-last, unlike recursive version:

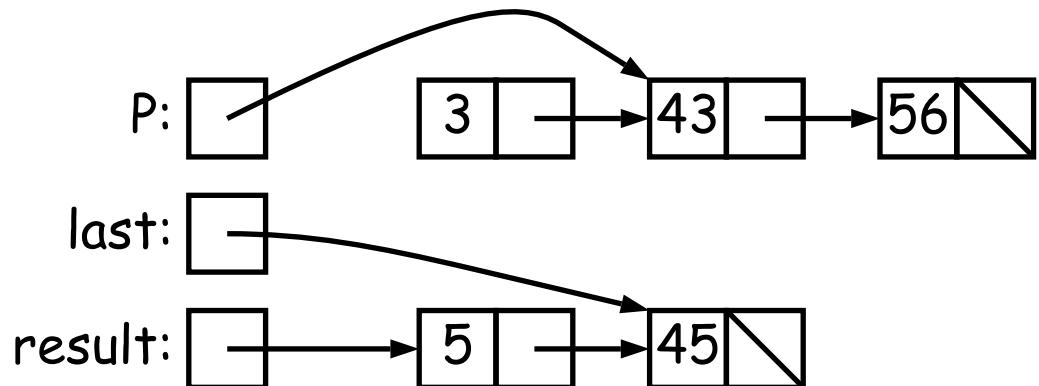
```
static IntList incrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    IntList result, last;  
  
    return result;  
}
```



An Iterative Version

An iterative incrList is tricky, because it is *not* tail recursive.
Easier to build things first-to-last, unlike recursive version:

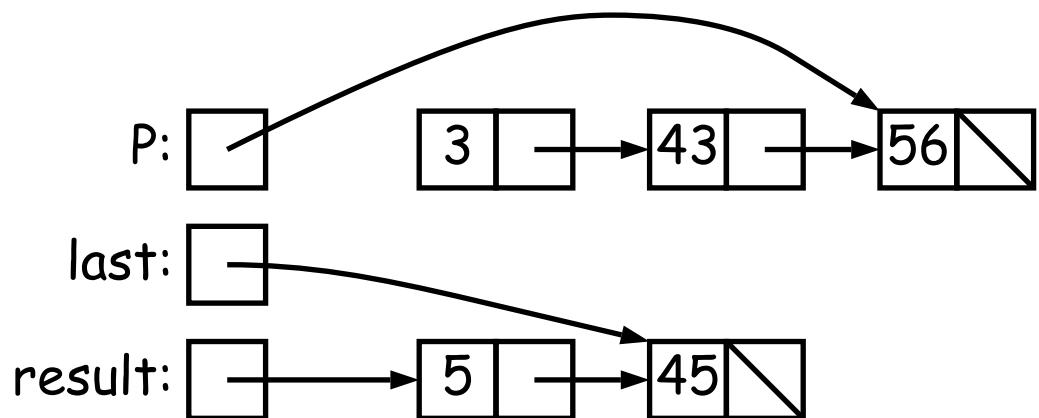
```
static IntList incrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    IntList result, last;  
  
    return result;  
}
```



An Iterative Version

An iterative incrList is tricky, because it is *not* tail recursive.
Easier to build things first-to-last, unlike recursive version:

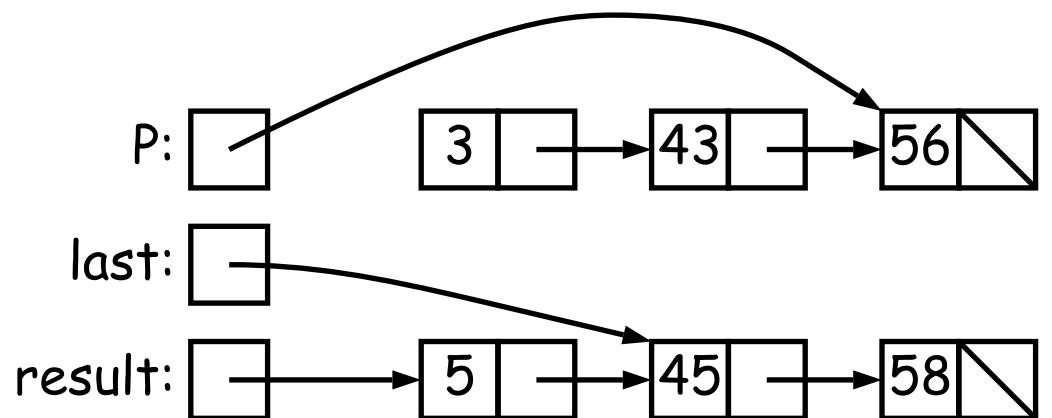
```
static IntList incrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    IntList result, last;  
  
    return result;  
}
```



An Iterative Version

An iterative incrList is tricky, because it is *not* tail recursive.
Easier to build things first-to-last, unlike recursive version:

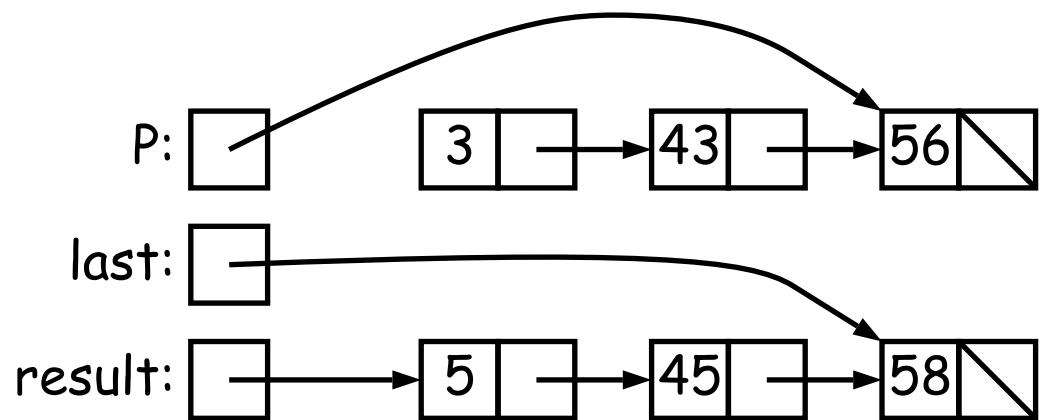
```
static IntList incrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    IntList result, last;  
  
    return result;  
}
```



An Iterative Version

An iterative incrList is tricky, because it is *not* tail recursive.
Easier to build things first-to-last, unlike recursive version:

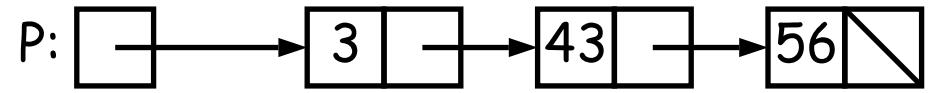
```
static IntList incrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    IntList result, last;  
  
    return result;  
}
```



An Iterative Version

An iterative incrList is tricky, because it is *not* tail recursive.
Easier to build things first-to-last, unlike recursive version:

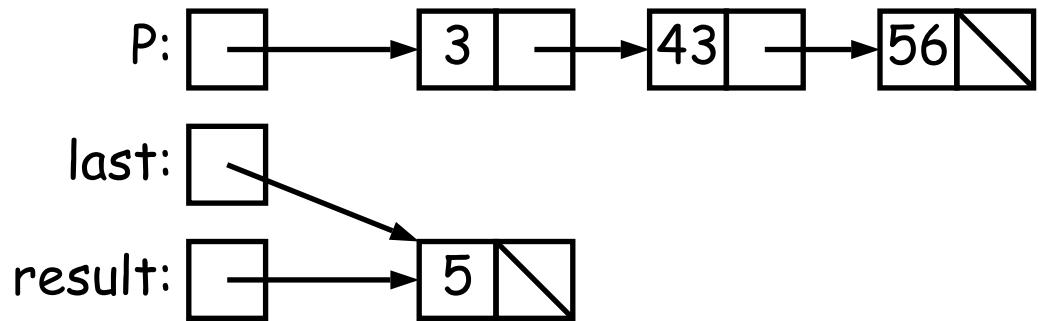
```
static IntList incrList(IntList P, int n) {  
    if (P == null)      <<<  
        return null;  
    IntList result, last;  
    result = last  
    = new IntList(P.head+n, null);  
    while (P.tail != null) {  
        P = P.tail;  
        last.tail  
        = new IntList(P.head+n, null);  
        last = last.tail;  
    }  
    return result;  
}
```



An Iterative Version

An iterative incrList is tricky, because it is *not* tail recursive.
Easier to build things first-to-last, unlike recursive version:

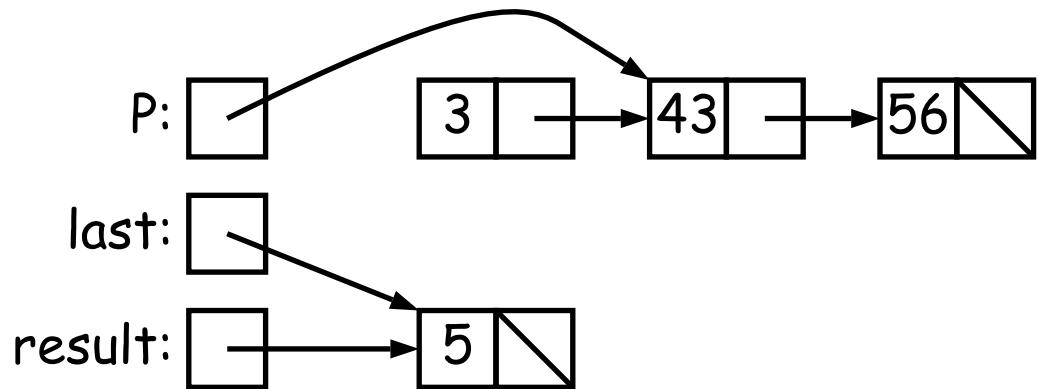
```
static IntList incrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    IntList result, last;  
    result = last      <<<  
    = new IntList(P.head+n, null);  
    while (P.tail != null) {  
        P = P.tail;  
        last.tail  
        = new IntList(P.head+n, null);  
        last = last.tail;  
    }  
    return result;  
}
```



An Iterative Version

An iterative incrList is tricky, because it is *not* tail recursive.
Easier to build things first-to-last, unlike recursive version:

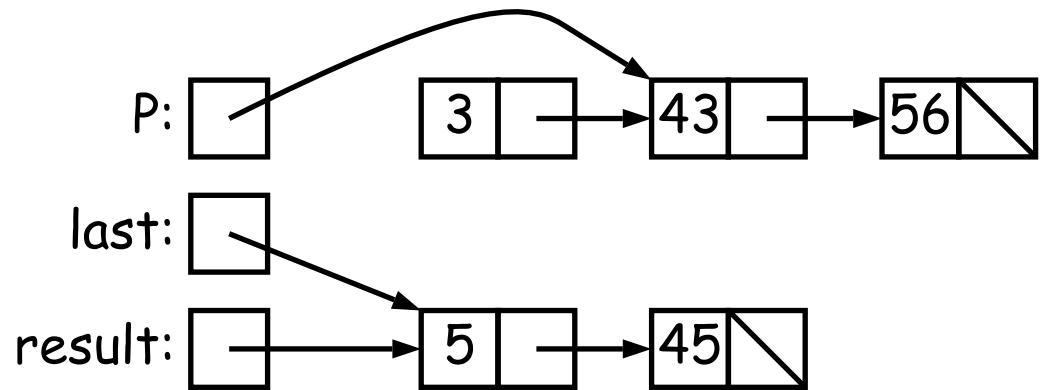
```
static IntList incrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    IntList result, last;  
    result = last  
        = new IntList(P.head+n, null);  
    while (P.tail != null) {  
        P = P.tail;      <<<  
        last.tail  
            = new IntList(P.head+n, null);  
        last = last.tail;  
    }  
    return result;  
}
```



An Iterative Version

An iterative incrList is tricky, because it is *not* tail recursive.
Easier to build things first-to-last, unlike recursive version:

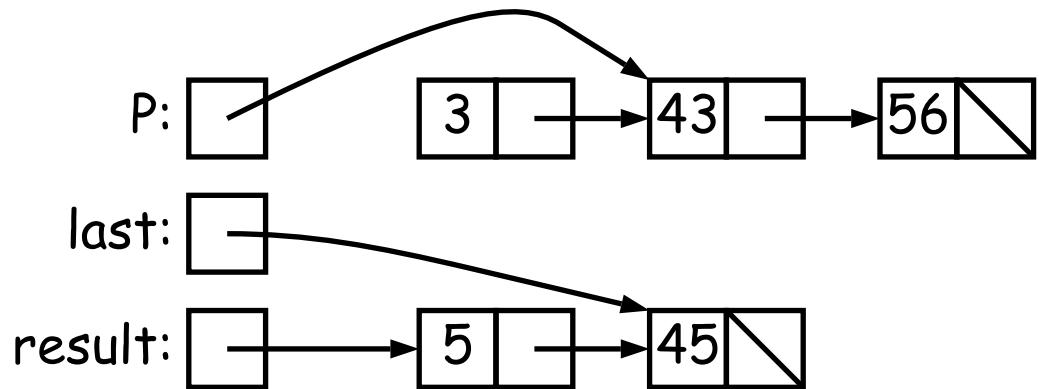
```
static IntList incrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    IntList result, last;  
    result = last  
        = new IntList(P.head+n, null);  
    while (P.tail != null) {  
        P = P.tail;  
        last.tail      <<<  
        = new IntList(P.head+n, null);  
        last = last.tail;  
    }  
    return result;  
}
```



An Iterative Version

An iterative incrList is tricky, because it is *not* tail recursive.
Easier to build things first-to-last, unlike recursive version:

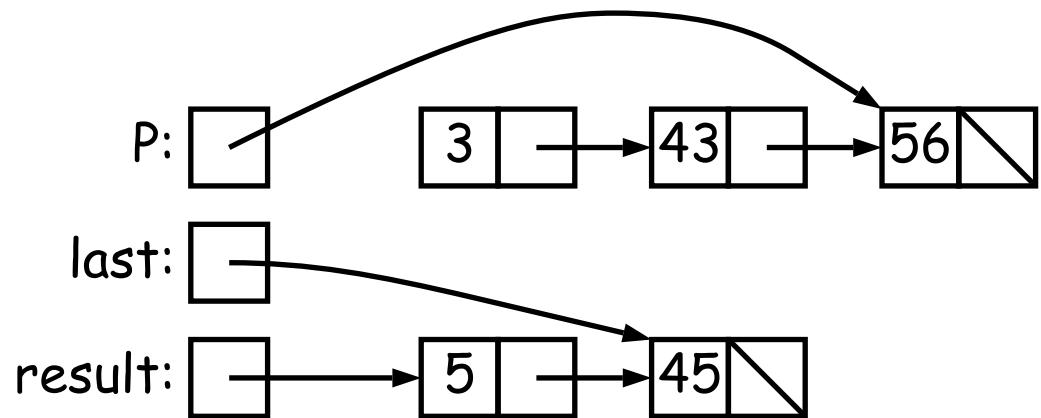
```
static IntList incrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    IntList result, last;  
    result = last  
        = new IntList(P.head+n, null);  
    while (P.tail != null) {  
        P = P.tail;  
        last.tail  
            = new IntList(P.head+n, null);  
        last = last.tail; <<<  
    }  
    return result;  
}
```



An Iterative Version

An iterative incrList is tricky, because it is *not* tail recursive.
Easier to build things first-to-last, unlike recursive version:

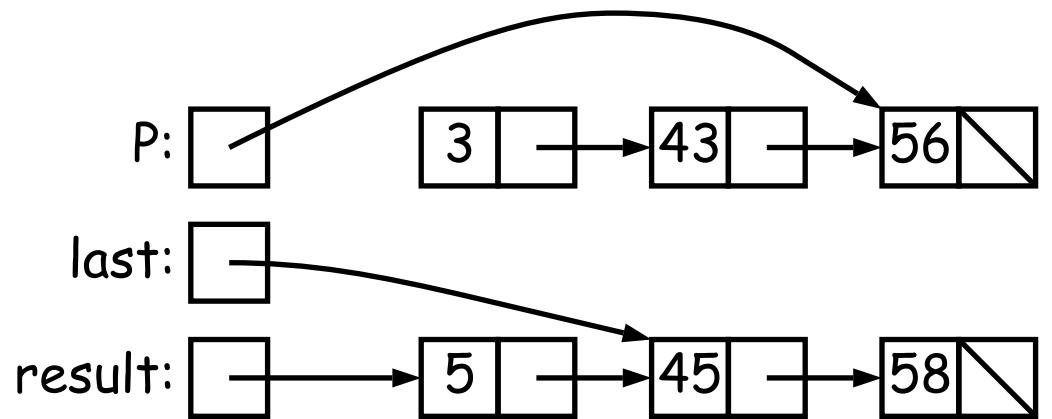
```
static IntList incrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    IntList result, last;  
    result = last  
        = new IntList(P.head+n, null);  
    while (P.tail != null) {  
        P = P.tail;      <<<  
        last.tail  
            = new IntList(P.head+n, null);  
        last = last.tail;  
    }  
    return result;  
}
```



An Iterative Version

An iterative incrList is tricky, because it is *not* tail recursive.
Easier to build things first-to-last, unlike recursive version:

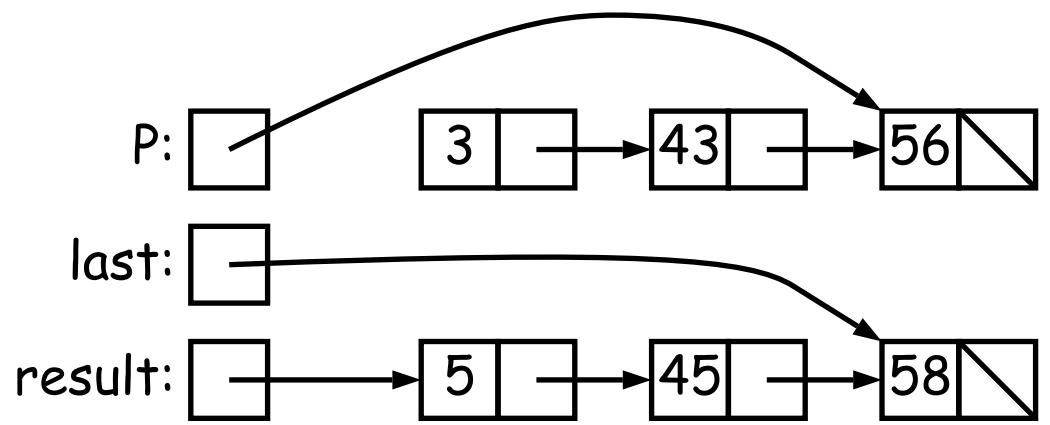
```
static IntList incrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    IntList result, last;  
    result = last  
        = new IntList(P.head+n, null);  
    while (P.tail != null) {  
        P = P.tail;  
        last.tail      <<<  
        = new IntList(P.head+n, null);  
        last = last.tail;  
    }  
    return result;  
}
```



An Iterative Version

An iterative incrList is tricky, because it is *not* tail recursive.
Easier to build things first-to-last, unlike recursive version:

```
static IntList incrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    IntList result, last;  
    result = last  
        = new IntList(P.head+n, null);  
    while (P.tail != null) {  
        P = P.tail;  
        last.tail  
            = new IntList(P.head+n, null);  
        last = last.tail; <<<  
    }  
    return result;  
}
```



CS61B Lecture #4: Simple Pointer Manipulation

Recreation Prove that for every acute angle $\alpha > 0$,

$$\tan \alpha + \cot \alpha \geq 2$$

Announcements

- Today: More pointer hacking.
- Handing in labs and homework: We'll be lenient about accepting late homework and labs for lab1, lab2, and hw0. Just get it done: part of the point is getting to understand the tools involved. We will **not** accept submissions by email.
- We will feel free to interpret the absence of a central repository for you or a lack of a lab1 submission from you as indicating that you intend to drop the course.

Small Test of Understanding

- In Java, the keyword `final` in a variable declaration means that the variable's value may not be changed after the variable is initialized.
- Is the following class valid?

```
public class Issue {  
  
    private final IntList aList = new IntList(0, null);  
  
    public void modify(int k) {  
        this.aList.head = k;  
    }  
}
```

Why or why not?

Small Test of Understanding

- In Java, the keyword `final` in a variable declaration means that the variable's value may not be changed after the variable is initialized.
- Is the following class valid?

```
public class Issue {  
  
    private final IntList aList = new IntList(0, null);  
  
    public void modify(int k) {  
        this.aList.head = k;  
    }  
}
```

Why or why not?

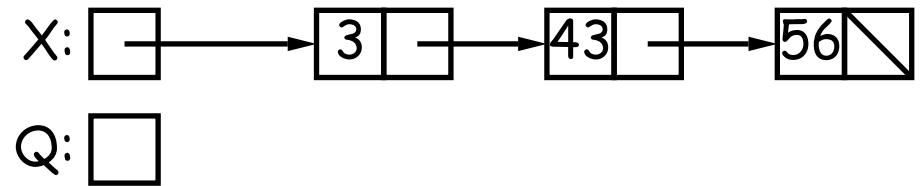
Answer: This is *valid*. Although `modify` changes the `head` variable of the object pointed to by `aList`, it does *not* modify the contents of `aList` itself (which is a pointer).

Destructive Incrementing

Destructive solutions may modify objects in the original list to save time or space:

```
/** Destructively add N to P's items.  
 * Return modified list. */  
static IntList dincrList(IntList P, int n) {  
}
```

```
X = IntList.list(3, 43, 56);  
/* IntList.list from HW #1 */  
Q = dincrList(X, 2);
```

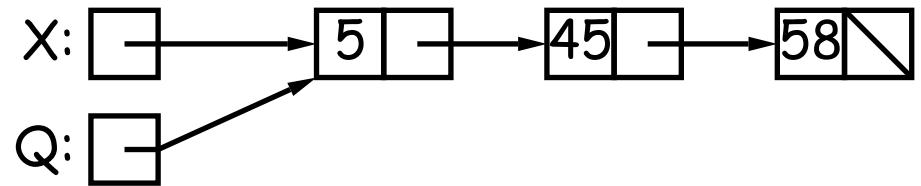


Destructive Incrementing

Destructive solutions may modify objects in the original list to save time or space:

```
/** Destructively add N to P's items.  
 * Return modified list. */  
static IntList dincrList(IntList P, int n) {  
}
```

```
X = IntList.list(3, 43, 56);  
/* IntList.list from HW #1 */  
Q = dincrList(X, 2);
```

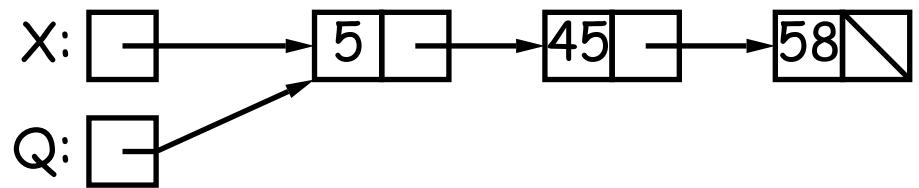


Destructive Incrementing

Destructive solutions may modify objects in the original list to save time or space:

```
/** Destructively add N to P's items.  
 * Return modified list. */  
static IntList dincrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    ?  
}
```

```
X = IntList.list(3, 43, 56);  
/* IntList.list from HW #1 */  
Q = dincrList(X, 2);
```

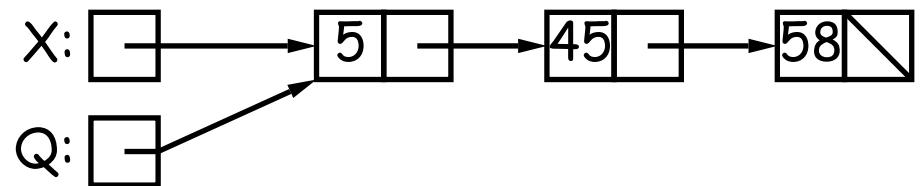


Destructive Incrementing

Destructive solutions may modify objects in the original list to save time or space:

```
/** Destructively add N to P's items.  
 * Return modified list. */  
static IntList dincrList(IntList P, int n) {  
    if (P == null)  
        return null;  
    else {  
        P.head += n;  
        P.tail = dincrList(P.tail, n);  
        return P;  
    }  
}
```

```
X = IntList.list(3, 43, 56);  
/* IntList.list from HW #1 */  
Q = dincrList(X, 2);
```

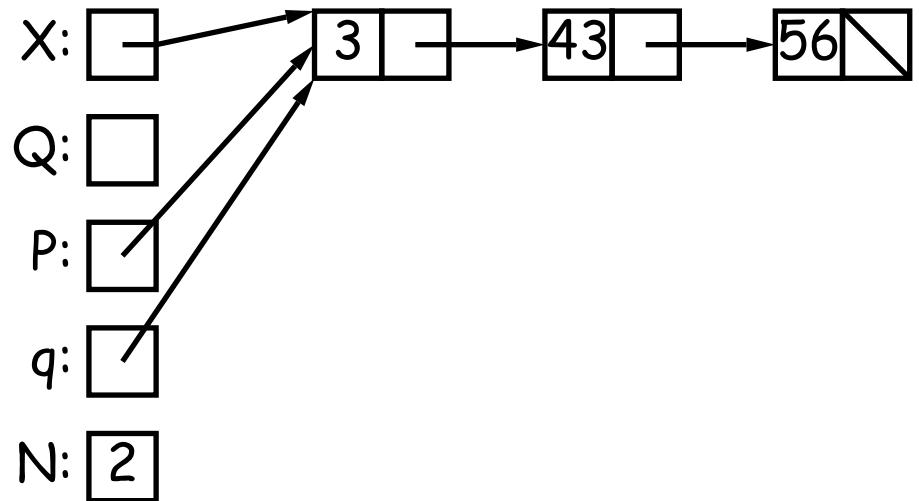


Iterative Destructive Incrementing

Let's try using a **while** loop:

```
/** Destructively add N to P's items.  
 * Return modified list. */  
static IntList dincrList(IntList P, int n) {  
}
```

```
X = IntList.list(3, 43, 56);  
/* IntList.list from HW #1 */  
Q = dincrList(X, 2);
```

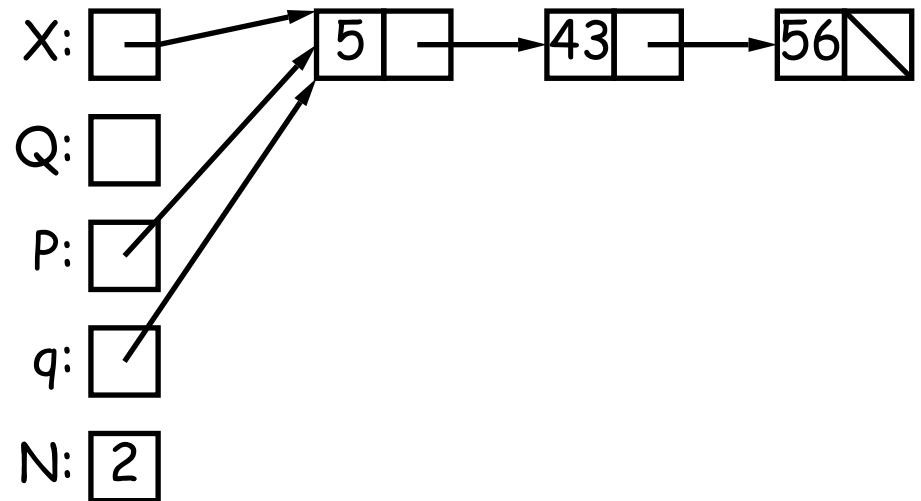


Iterative Destructive Incrementing

Let's try using a **while** loop:

```
/** Destructively add N to P's items.  
 * Return modified list. */  
static IntList dincrList(IntList P, int n) {  
}
```

```
X = IntList.list(3, 43, 56);  
/* IntList.list from HW #1 */  
Q = dincrList(X, 2);
```

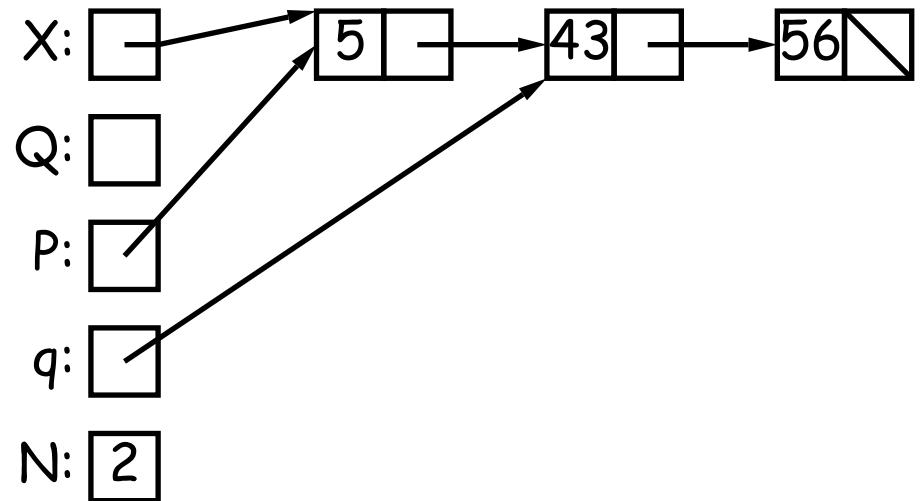


Iterative Destructive Incrementing

Let's try using a **while** loop:

```
/** Destructively add N to P's items.  
 * Return modified list. */  
static IntList dincrList(IntList P, int n) {  
}
```

```
X = IntList.list(3, 43, 56);  
/* IntList.list from HW #1 */  
Q = dincrList(X, 2);
```

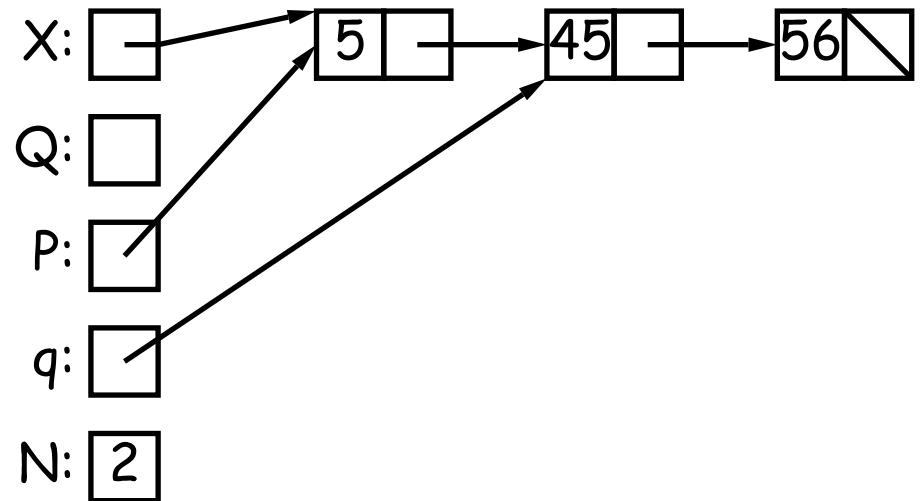


Iterative Destructive Incrementing

Let's try using a **while** loop:

```
/** Destructively add N to P's items.  
 * Return modified list. */  
static IntList dincrList(IntList P, int n) {  
}
```

```
X = IntList.list(3, 43, 56);  
/* IntList.list from HW #1 */  
Q = dincrList(X, 2);
```

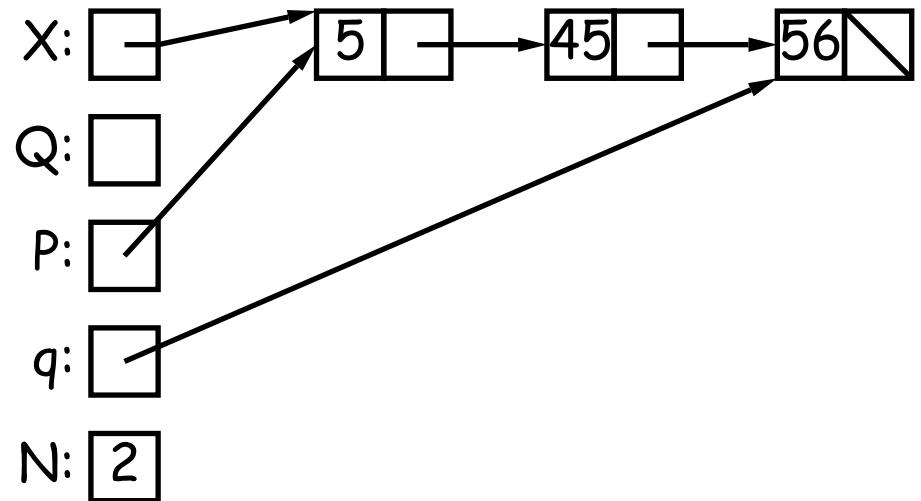


Iterative Destructive Incrementing

Let's try using a **while** loop:

```
/** Destructively add N to P's items.  
 * Return modified list. */  
static IntList dincrList(IntList P, int n) {  
}
```

```
X = IntList.list(3, 43, 56);  
/* IntList.list from HW #1 */  
Q = dincrList(X, 2);
```

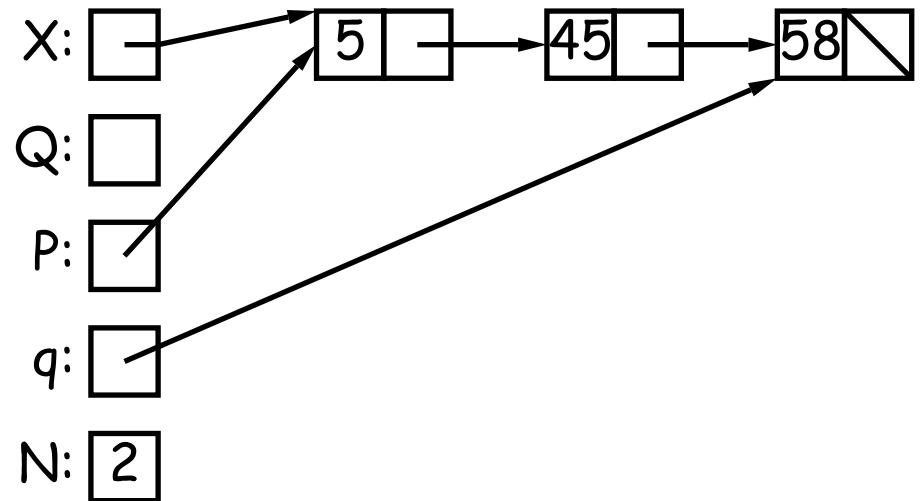


Iterative Destructive Incrementing

Let's try using a **while** loop:

```
/** Destructively add N to P's items.  
 * Return modified list. */  
static IntList dincrList(IntList P, int n) {  
}
```

```
X = IntList.list(3, 43, 56);  
/* IntList.list from HW #1 */  
Q = dincrList(X, 2);
```

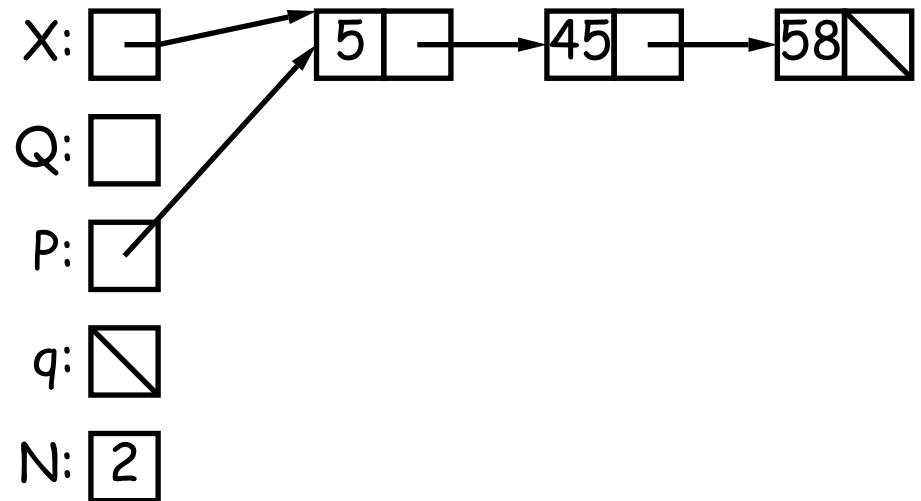


Iterative Destructive Incrementing

Let's try using a **while** loop:

```
/** Destructively add N to P's items.  
 * Return modified list. */  
static IntList dincrList(IntList P, int n) {  
}
```

```
X = IntList.list(3, 43, 56);  
/* IntList.list from HW #1 */  
Q = dincrList(X, 2);
```

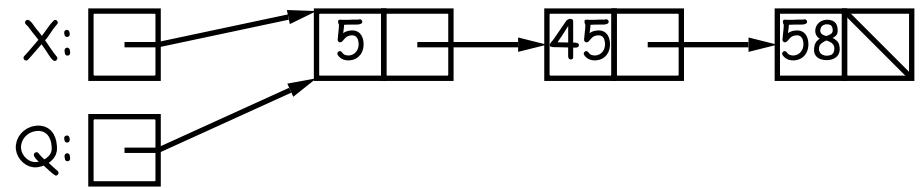


Iterative Destructive Incrementing

Let's try using a `while` loop:

```
/** Destructively add N to P's items.  
 * Return modified list. */  
static IntList dincrList(IntList P, int n) {  
    IntList q;  
    q = P;  
    while (?) {  
    }  
    return ?;  
}
```

```
X = IntList.list(3, 43, 56);  
/* IntList.list from HW #1 */  
Q = dincrList(X, 2);
```

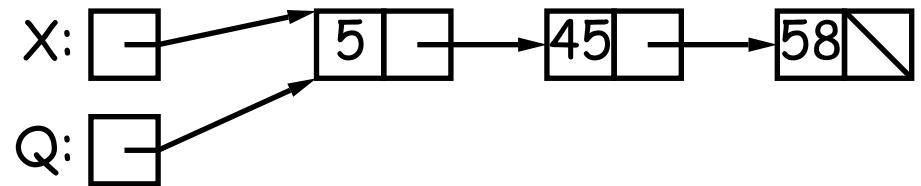


Iterative Destructive Incrementing

Let's try using a `while` loop:

```
/** Destructively add N to P's items.  
 * Return modified list. */  
static IntList dincrList(IntList P, int n) {  
    IntList q;  
    q = P;  
    while (q != null) {  
        q.head += n;  
        q = q.tail;  
    }  
    return P;  
}
```

```
X = IntList.list(3, 43, 56);  
/* IntList.list from HW #1 */  
Q = dincrList(X, 2);
```

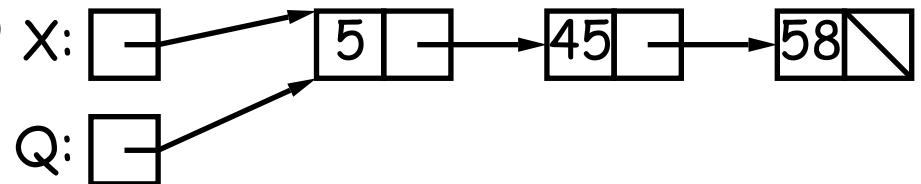


Iterative Destructive Incrementing

Let's try using a **while** loop:

```
/** Destructively add N to P's items.  
 * Return modified list. */  
static IntList dincrList(IntList P, int n) {  
    // 'for' can do more than count!  
    for (IntList q = P; q != null; q = q.tail)  
        q.head += n;  
    return P;  
}
```

```
X = IntList.list(3, 43, 56);  
/* IntList.list from HW #1 */  
Q = dincrList(X, 2);
```



Another Example: Non-destructive List Deletion

If L is the list [2, 1, 2, 9, 2], we want removeAll(L, 2) to be the new list [1, 9].

```
/** The list resulting from removing all instances of X from L
 * non-destructively. */
static IntList removeAll(IntList L, int x) {
    if (L == null)
        return /*( null with all x's removed )*/;
    else if (L.head == x)
        return /*( L with all x's removed (L!=null, L.head==x) )*/;
    else
        return /*( L with all x's removed (L!=null, L.head!=x) )*/;
}
```

Another Example: Non-destructive List Deletion

If L is the list [2, 1, 2, 9, 2], we want removeAll(L, 2) to be the new list [1, 9].

```
/** The list resulting from removing all instances of X from L
 * non-destructively. */
static IntList removeAll(IntList L, int x) {
    if (L == null)
        return null;
    else if (L.head == x)
        return /*( L with all x's removed (L!=null, L.head==x) );*/
    else
        return /*( L with all x's removed (L!=null, L.head!=x) );*/
}
```

Another Example: Non-destructive List Deletion

If L is the list [2, 1, 2, 9, 2], we want removeAll(L, 2) to be the new list [1, 9].

```
/** The list resulting from removing all instances of X from L
 * non-destructively. */
static IntList removeAll(IntList L, int x) {
    if (L == null)
        return null;
    else if (L.head == x)
        return removeAll(L.tail, x);
    else
        return /*( L with all x's removed (L!=null, L.head!=x) );*/
}
```

Another Example: Non-destructive List Deletion

If L is the list [2, 1, 2, 9, 2], we want removeAll(L, 2) to be the new list [1, 9].

```
/** The list resulting from removing all instances of X from L
 * non-destructively. */
static IntList removeAll(IntList L, int x) {
    if (L == null)
        return null;
    else if (L.head == x)
        return removeAll(L.tail, x);
    else
        return new IntList(L.head, removeAll(L.tail, x));
}
```

Iterative Non-destructive List Deletion

Same as before, but use front-to-back iteration rather than recursion.

```
/** The list resulting from removing all instances
 * of X from L non-destructively. */
static IntList removeAll(IntList L, int x) {
    IntList result, last;
    result = last = null;

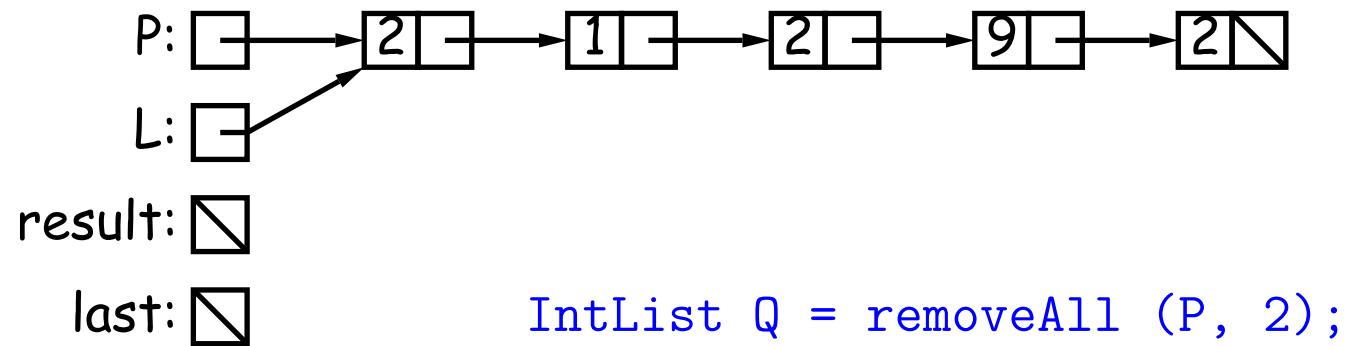
    ?

    return result;
}
```

Iterative Non-destructive List Deletion

Same as before, but use front-to-back iteration rather than recursion.

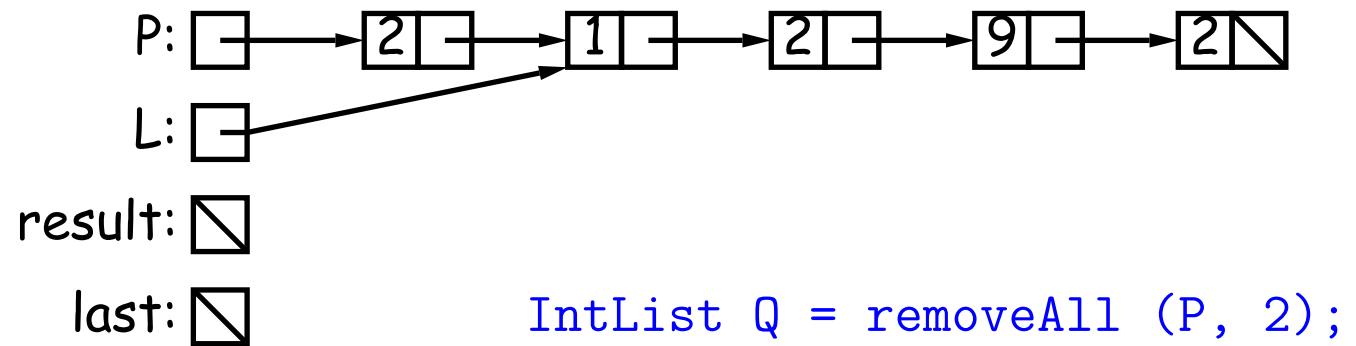
```
/** The list resulting from removing all instances
 * of X from L non-destructively. */
static IntList removeAll(IntList L, int x) {
    IntList result, last;
    result = last = null;
    ?
    return result;
}
```



Iterative Non-destructive List Deletion

Same as before, but use front-to-back iteration rather than recursion.

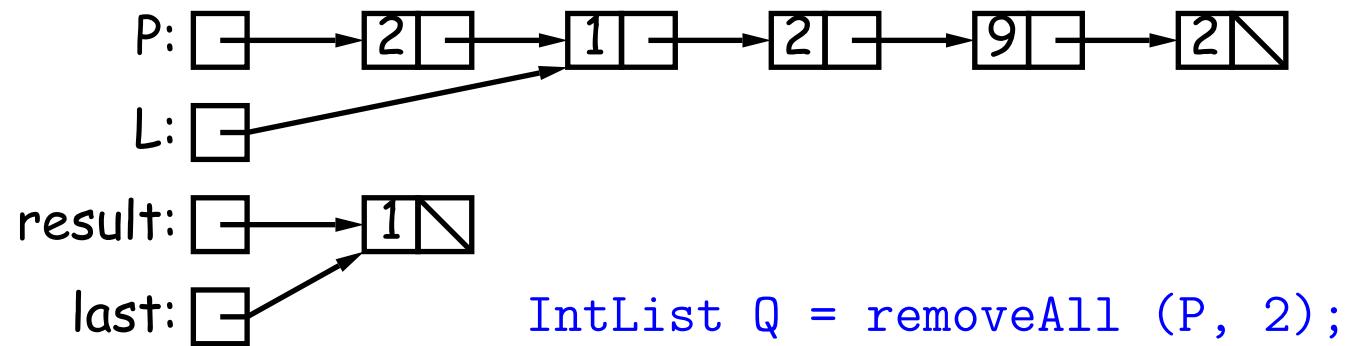
```
/** The list resulting from removing all instances
 * of X from L non-destructively. */
static IntList removeAll(IntList L, int x) {
    IntList result, last;
    result = last = null;
    ?
    return result;
}
```



Iterative Non-destructive List Deletion

Same as before, but use front-to-back iteration rather than recursion.

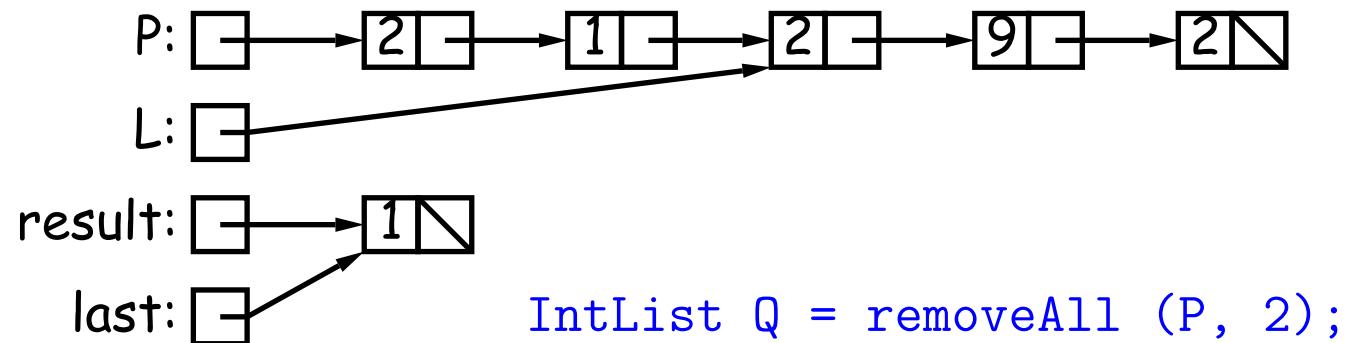
```
/** The list resulting from removing all instances
 * of X from L non-destructively. */
static IntList removeAll(IntList L, int x) {
    IntList result, last;
    result = last = null;
    ?
    return result;
}
```



Iterative Non-destructive List Deletion

Same as before, but use front-to-back iteration rather than recursion.

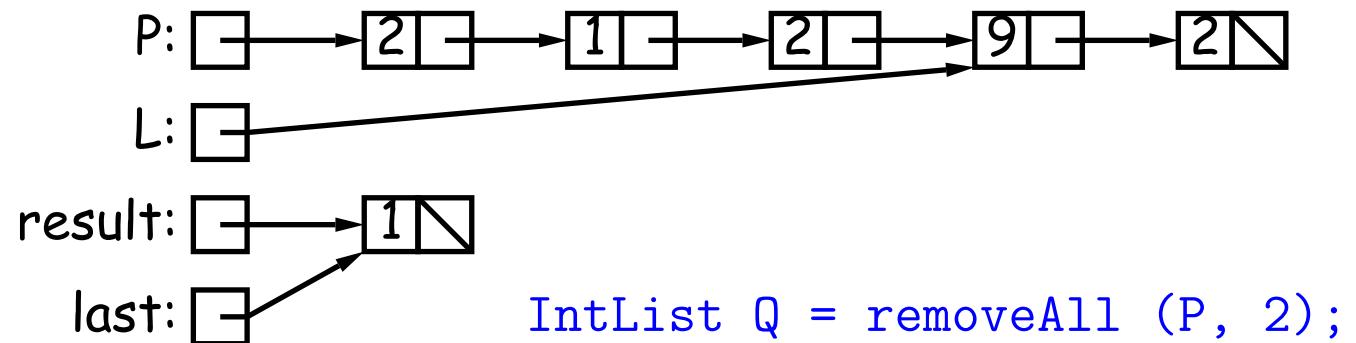
```
/** The list resulting from removing all instances
 * of X from L non-destructively. */
static IntList removeAll(IntList L, int x) {
    IntList result, last;
    result = last = null;
    ?
    return result;
}
```



Iterative Non-destructive List Deletion

Same as before, but use front-to-back iteration rather than recursion.

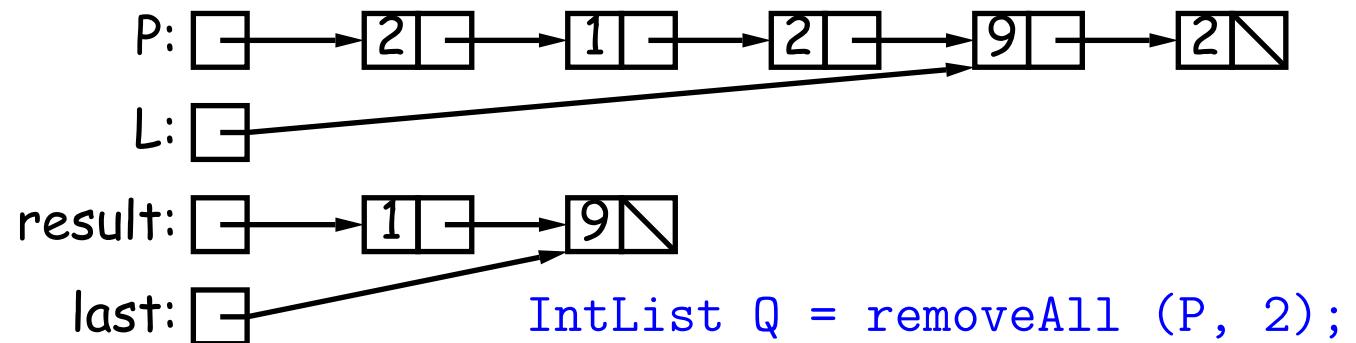
```
/** The list resulting from removing all instances
 * of X from L non-destructively. */
static IntList removeAll(IntList L, int x) {
    IntList result, last;
    result = last = null;
    ?
    return result;
}
```



Iterative Non-destructive List Deletion

Same as before, but use front-to-back iteration rather than recursion.

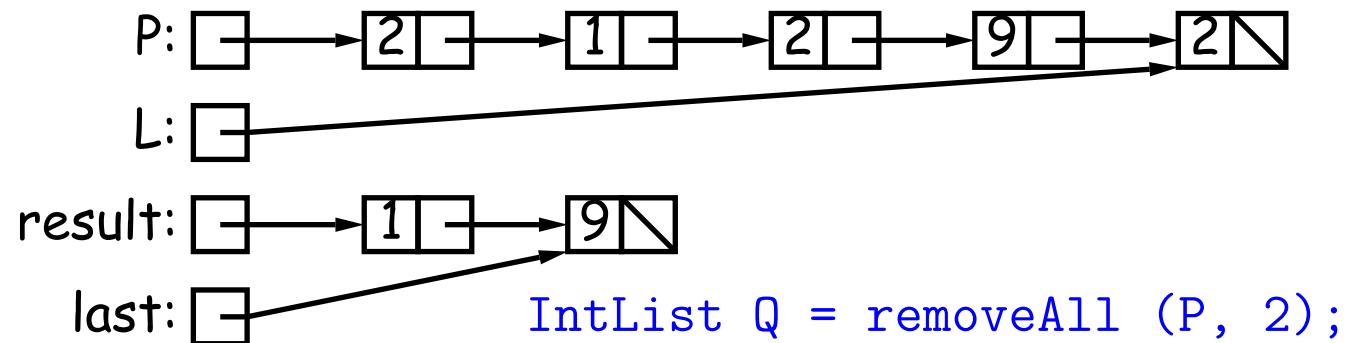
```
/** The list resulting from removing all instances
 * of X from L non-destructively. */
static IntList removeAll(IntList L, int x) {
    IntList result, last;
    result = last = null;
    ?
    return result;
}
```



Iterative Non-destructive List Deletion

Same as before, but use front-to-back iteration rather than recursion.

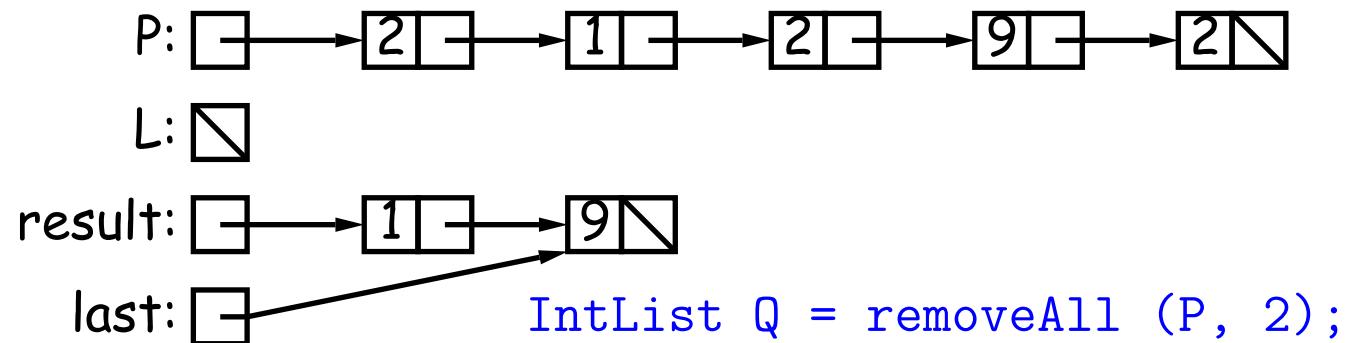
```
/** The list resulting from removing all instances
 * of X from L non-destructively. */
static IntList removeAll(IntList L, int x) {
    IntList result, last;
    result = last = null;
    ?
    return result;
}
```



Iterative Non-destructive List Deletion

Same as before, but use front-to-back iteration rather than recursion.

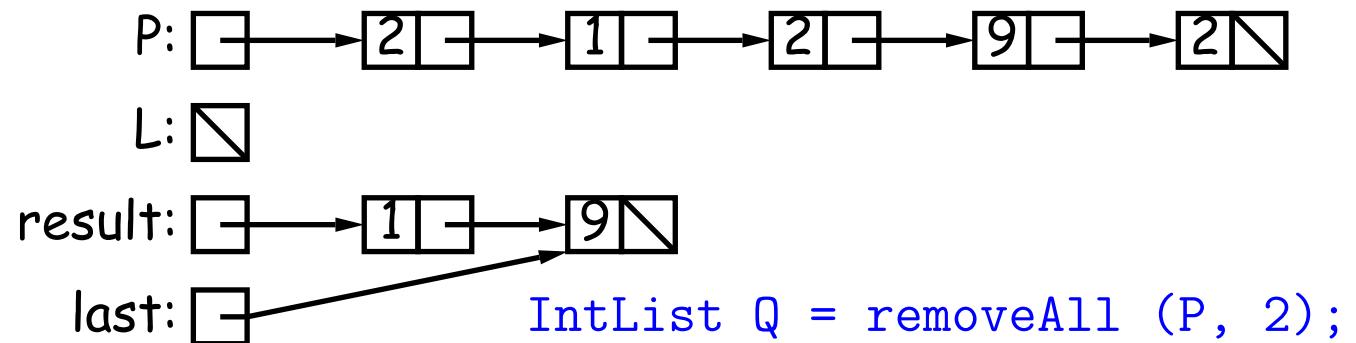
```
/** The list resulting from removing all instances
 * of X from L non-destructively. */
static IntList removeAll(IntList L, int x) {
    IntList result, last;
    result = last = null;
    ?
    return result;
}
```



Iterative Non-destructive List Deletion

Same as before, but use front-to-back iteration rather than recursion.

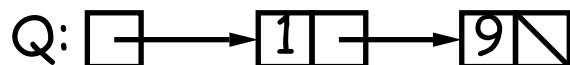
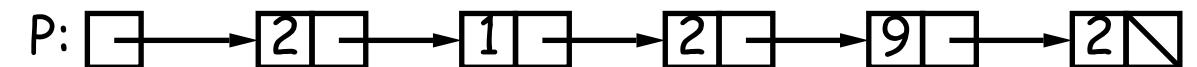
```
/** The list resulting from removing all instances
 * of X from L non-destructively. */
static IntList removeAll(IntList L, int x) {
    IntList result, last;
    result = last = null;
    for ( ; L != null; L = L.tail) {
        if (x == L.head)
            continue;
        else if (last == null)
            result = last = new IntList(L.head, null);
        else
            last = last.tail = new IntList(L.head, null);
    }
    return result;
}
```



Iterative Non-destructive List Deletion

Same as before, but use front-to-back iteration rather than recursion.

```
/** The list resulting from removing all instances
 * of X from L non-destructively. */
static IntList removeAll(IntList L, int x) {
    IntList result, last;
    result = last = null;
    for ( ; L != null; L = L.tail) {
        if (x == L.head)
            continue;
        else if (last == null)
            result = last = new IntList(L.head, null);
        else
            last = last.tail = new IntList(L.head, null);
    }
    return result;
}
```



IntList Q = removeAll (P, 2);

Destructive Deletion

→ : Original

..... : after $Q = \text{dremoveAll}(Q, 1)$

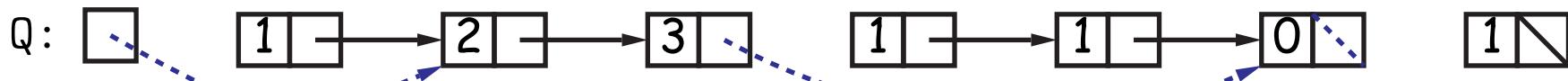


```
/** The list resulting from removing all instances of X from L.
 * The original list may be destroyed. */
static IntList dremoveAll(IntList L, int x) {
    if (L == null)
        return /*( null with all x's removed )*/;
    else if (L.head == x)
        return /*( L with all x's removed (L != null) )*/;
    else {
        /*{ Remove all x's from L's tail. }*/;
        return L;
    }
}
```

Destructive Deletion

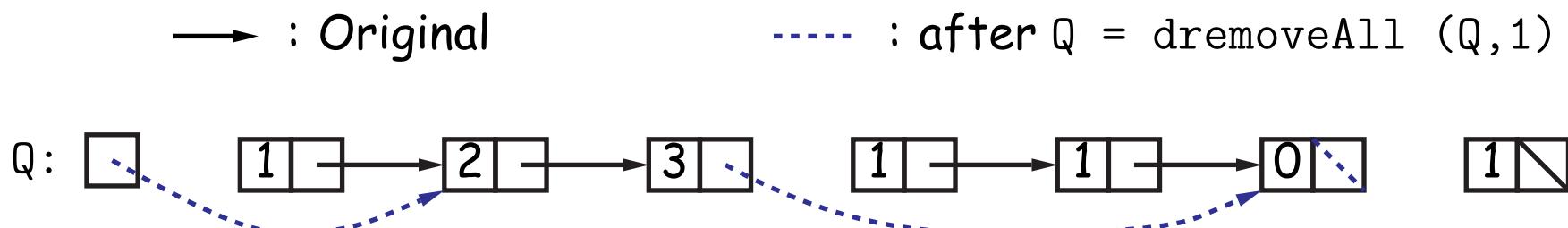
→ : Original

..... : after $Q = \text{dremoveAll}(Q, 1)$



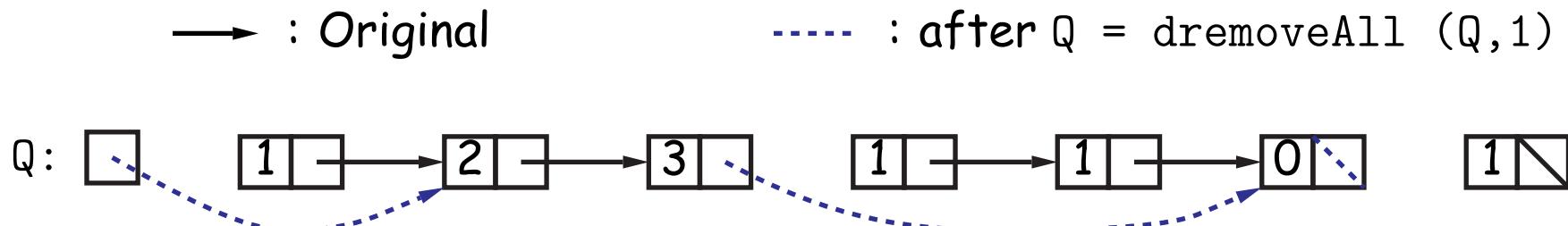
```
/** The list resulting from removing all instances of X from L.
 * The original list may be destroyed. */
static IntList dremoveAll(IntList L, int x) {
    if (L == null)
        return /*( null with all x's removed )*/;
    else if (L.head == x)
        return /*( L with all x's removed (L != null) )*/;
    else {
        /*{ Remove all x's from L's tail. }*/;
        return L;
    }
}
```

Destructive Deletion



```
/** The list resulting from removing all instances of X from L.
 * The original list may be destroyed. */
static IntList dremoveAll(IntList L, int x) {
    if (L == null)
        return /*( null with all x's removed )*/;
    else if (L.head == x)
        return /*( L with all x's removed (L != null) )*/;
    else {
        /*{ Remove all x's from L's tail. }*/;
        return L;
    }
}
```

Destructive Deletion

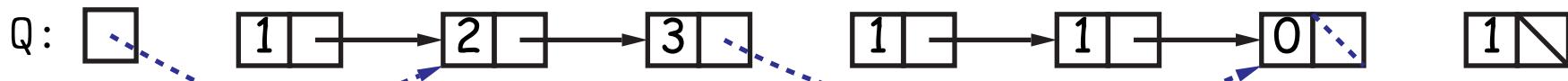


```
/** The list resulting from removing all instances of X from L.
 * The original list may be destroyed. */
static IntList dremoveAll(IntList L, int x) {
    if (L == null)
        return /*( null with all x's removed )*/;
    else if (L.head == x)
        return /*( L with all x's removed (L != null) )*/;
    else {
        /*{ Remove all x's from L's tail. }*/;
        return L;
    }
}
```

Destructive Deletion

→ : Original

..... : after $Q = \text{dremoveAll}(Q, 1)$

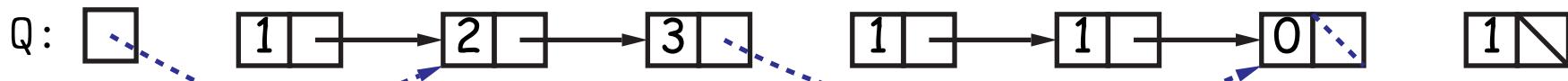


```
/** The list resulting from removing all instances of X from L.
 * The original list may be destroyed. */
static IntList dremoveAll(IntList L, int x) {
    if (L == null)
        return null;
    else if (L.head == x)
        return /*( L with all x's removed (L != null) );*/
    else {
        /*{ Remove all x's from L's tail. }*/;
        return L;
    }
}
```

Destructive Deletion

→ : Original

..... : after $Q = \text{dremoveAll}(Q, 1)$



```
/** The list resulting from removing all instances of X from L.
 * The original list may be destroyed. */
static IntList dremoveAll(IntList L, int x) {
    if (L == null)
        return
    else if (L.head == x)
        return dremoveAll(L.tail, x);
    else {
        /*{ Remove all x's from L's tail. }*/;
        return L;
    }
}
```

Destructive Deletion

→ : Original

..... : after $Q = \text{dremoveAll}(Q, 1)$



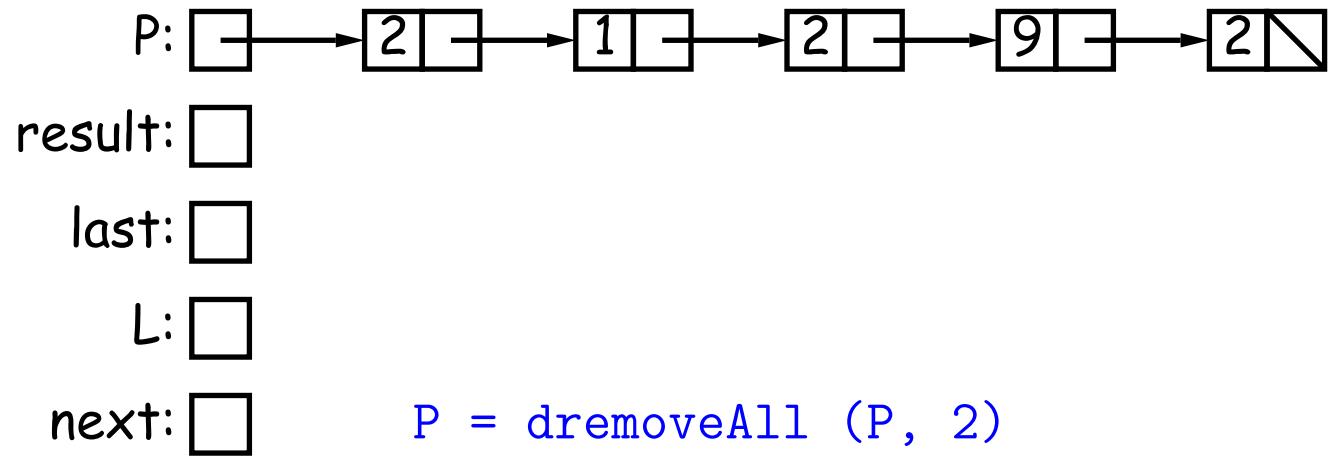
```
/** The list resulting from removing all instances of X from L.  
 * The original list may be destroyed. */  
static IntList dremoveAll(IntList L, int x) {  
    if (L == null)  
        return  
    else if (L.head == x)  
        return dremoveAll(L.tail, x);  
    else {  
        L.tail = dremoveAll(L.tail, x);  
        return L;  
    }  
}
```

Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
        }
    return result;
}
```

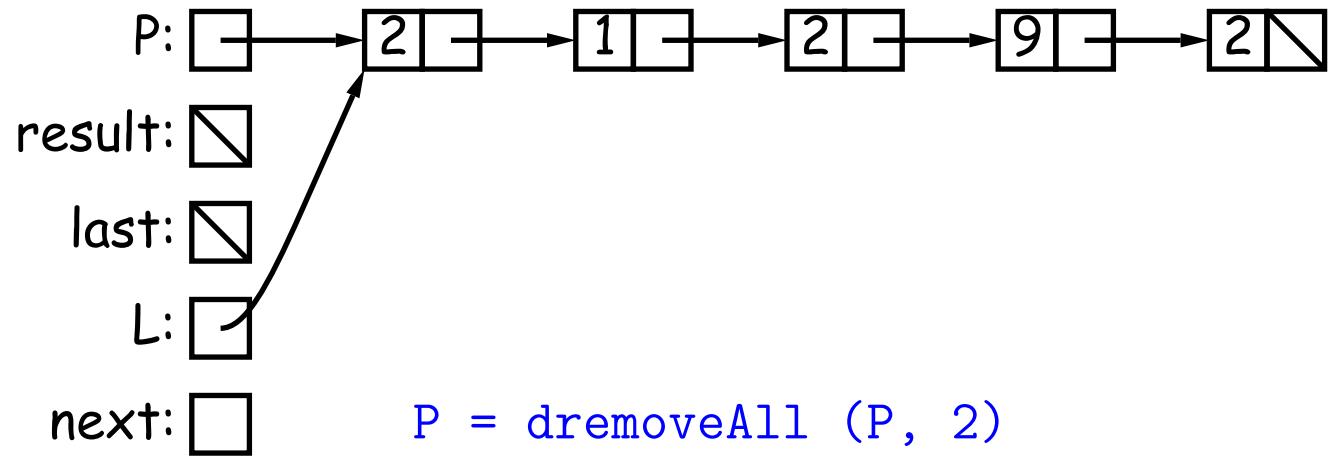
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



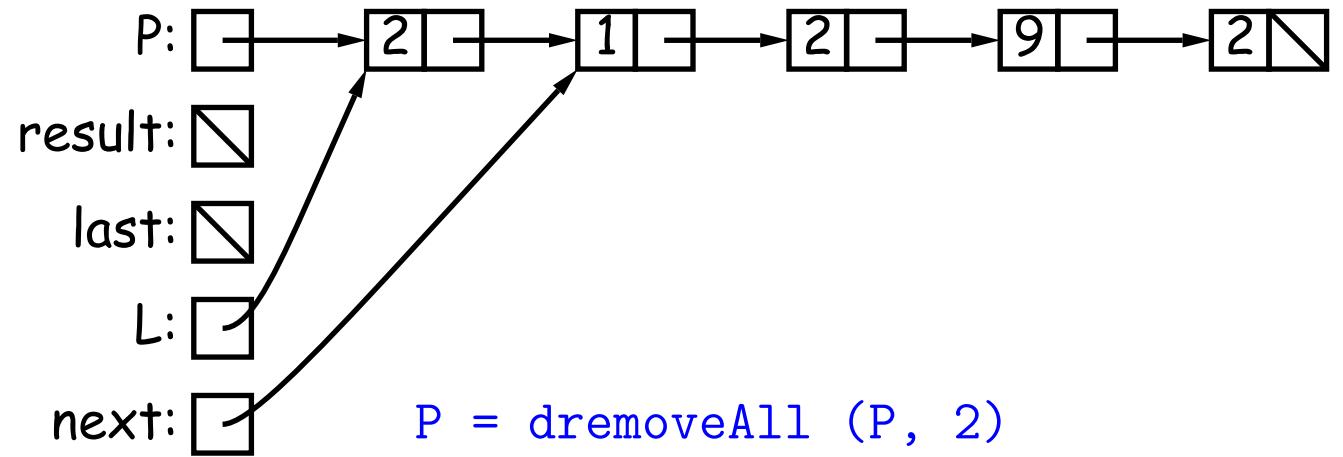
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



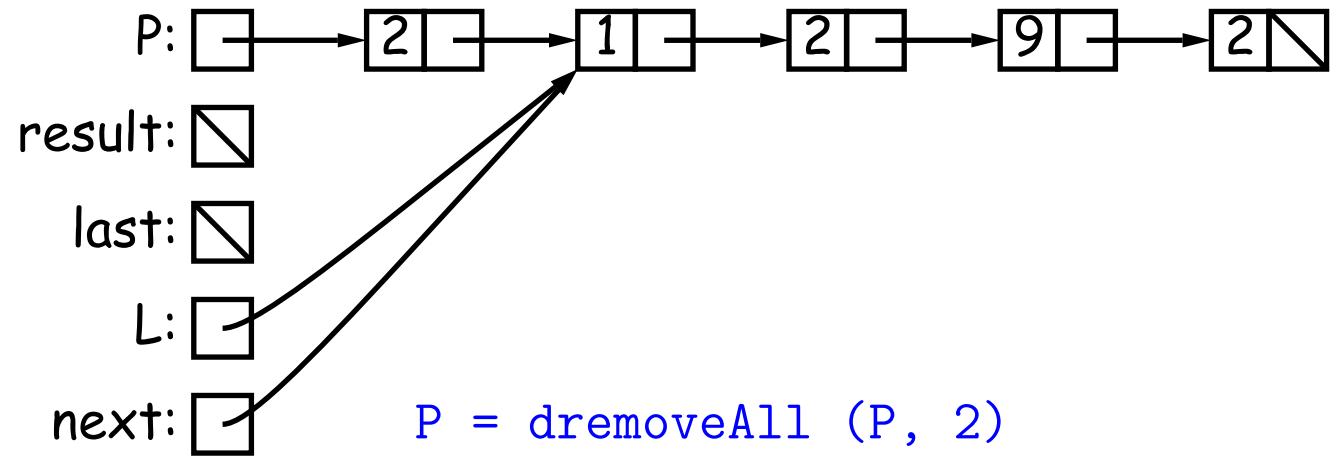
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



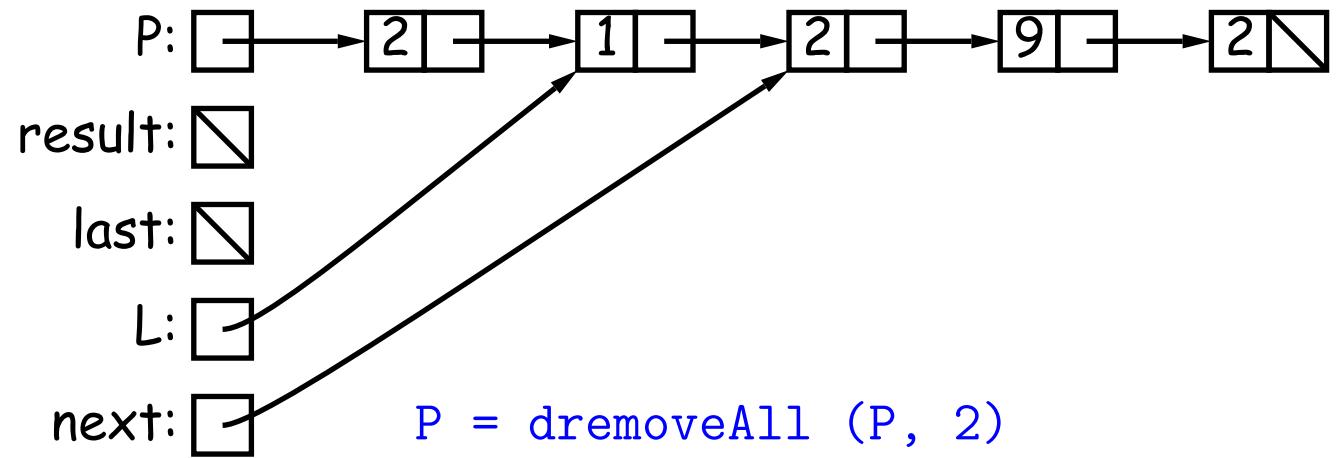
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



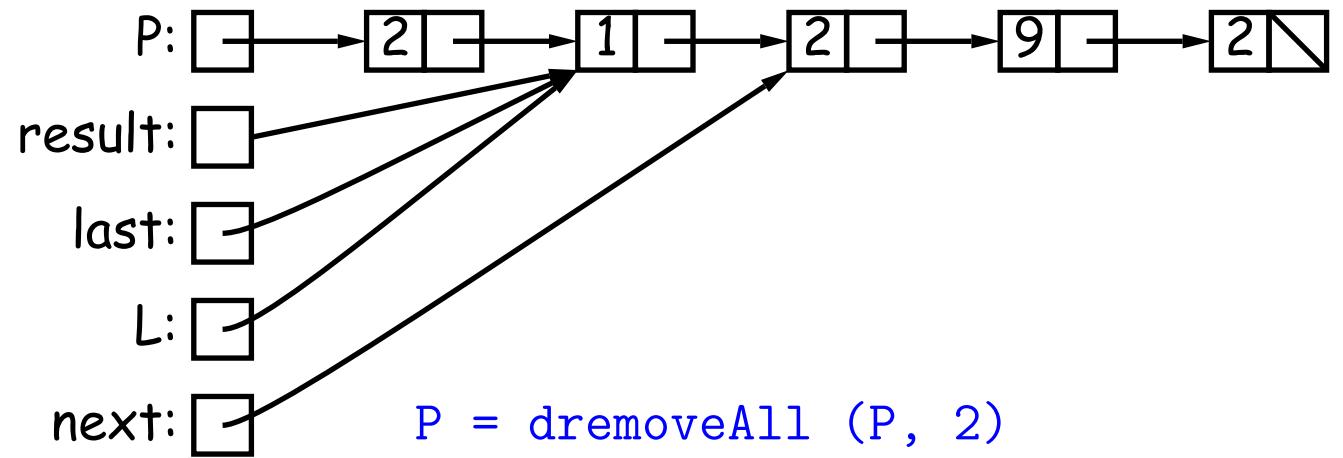
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



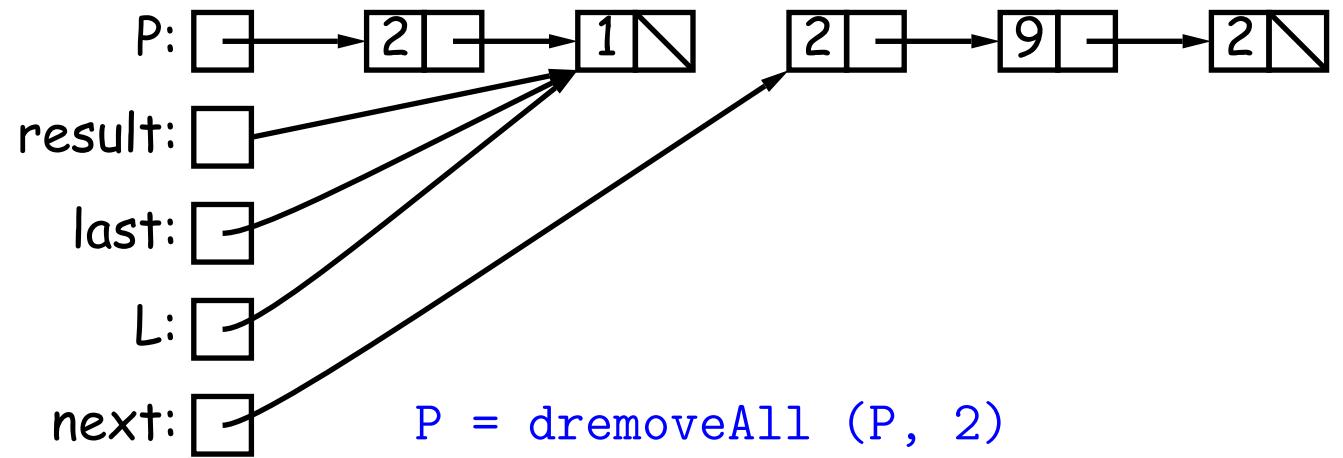
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



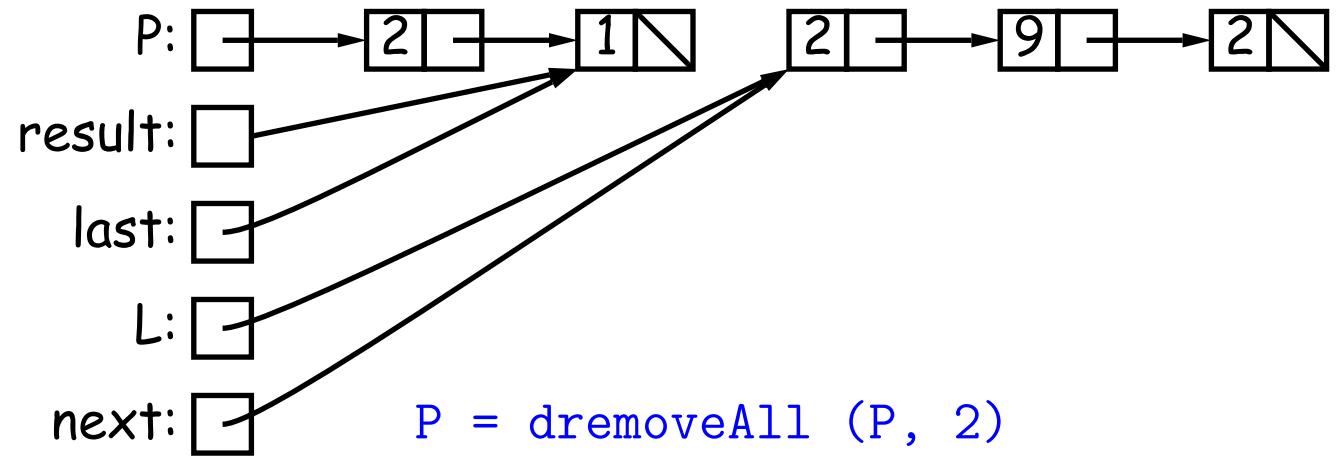
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



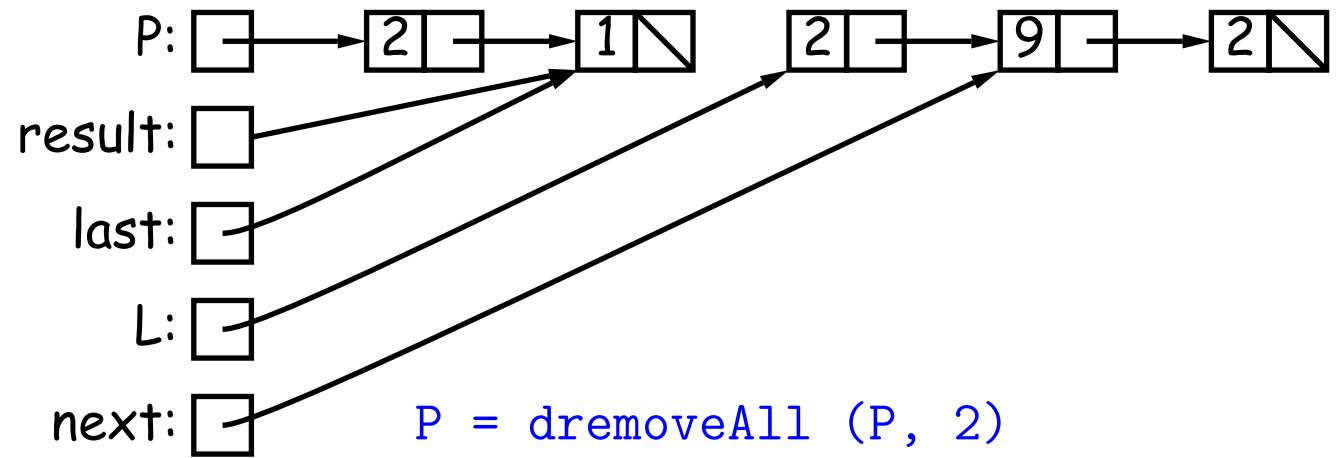
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



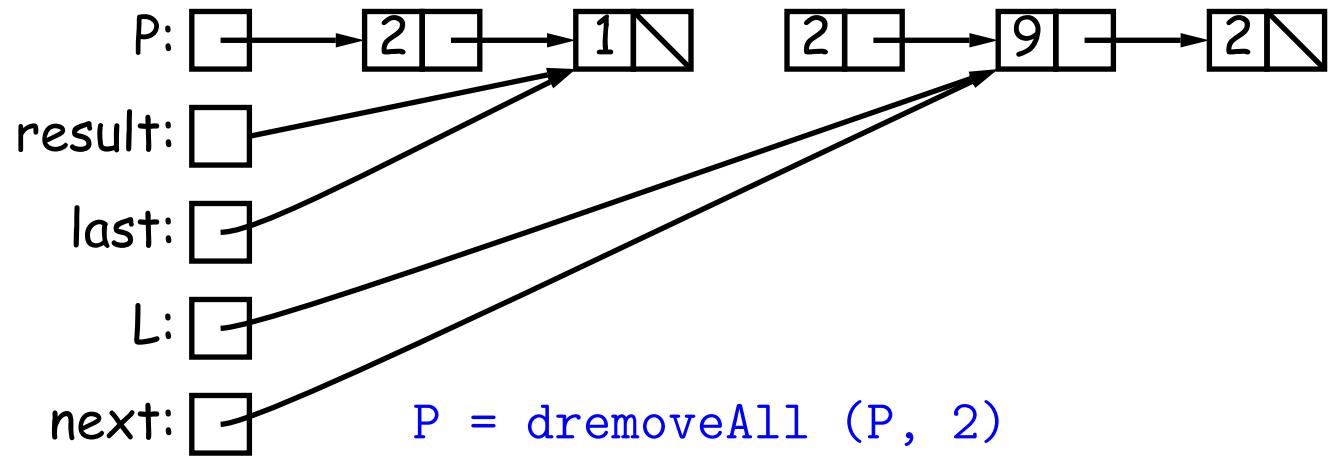
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



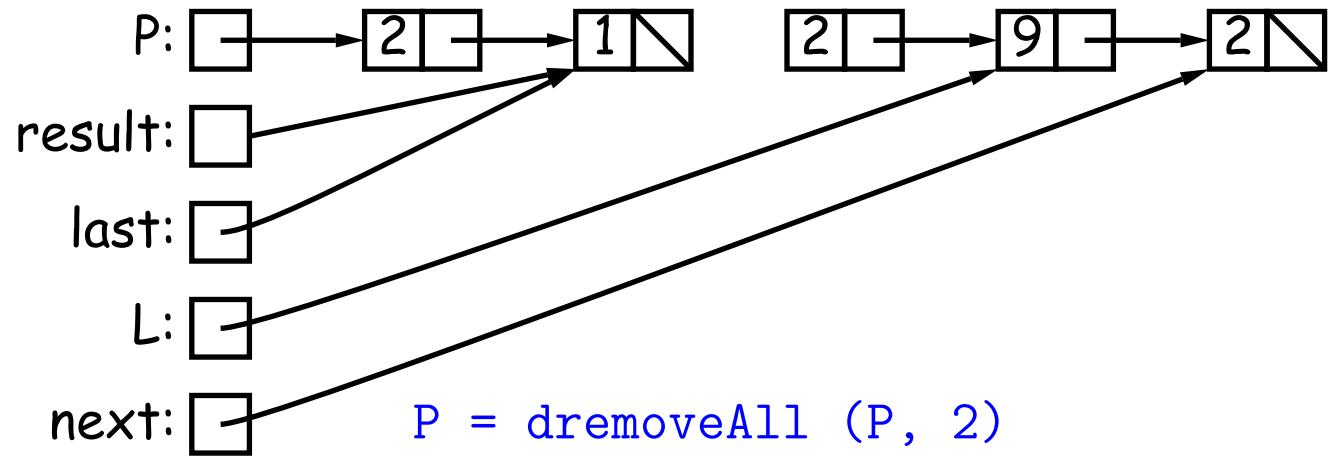
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



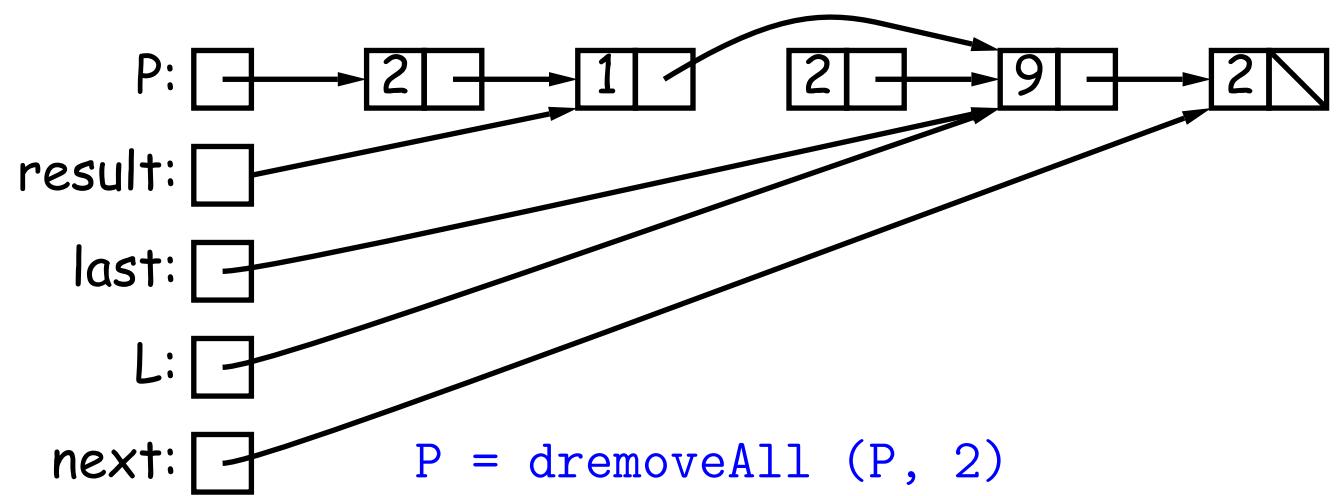
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



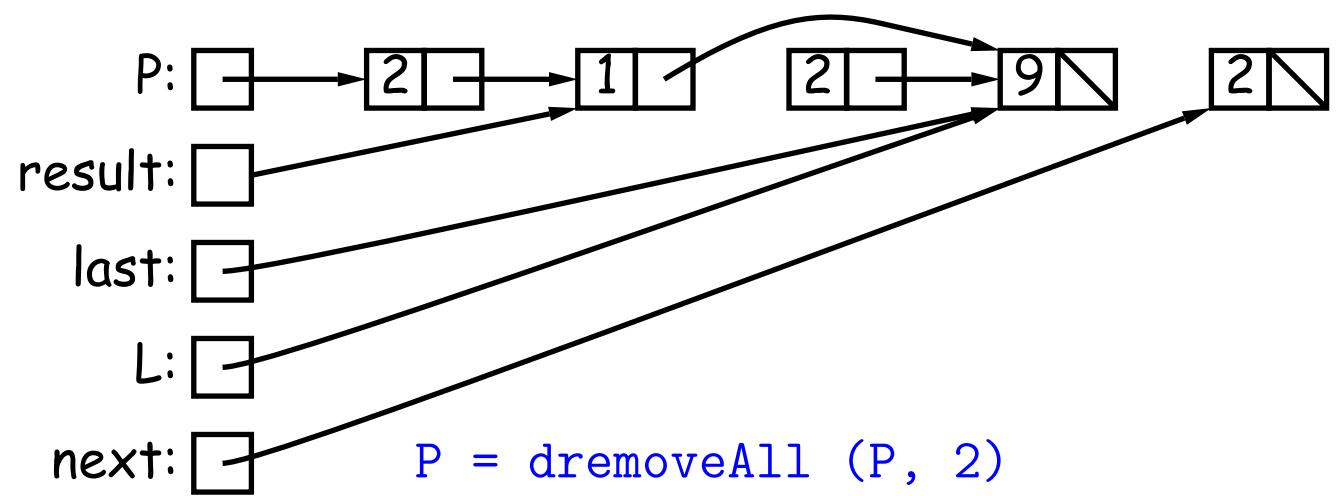
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



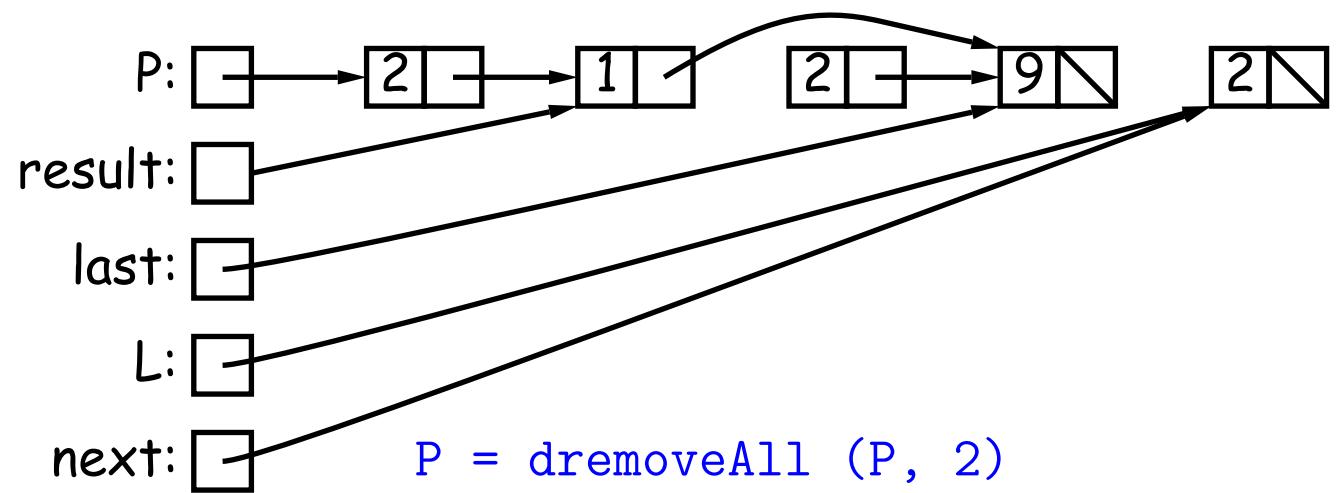
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



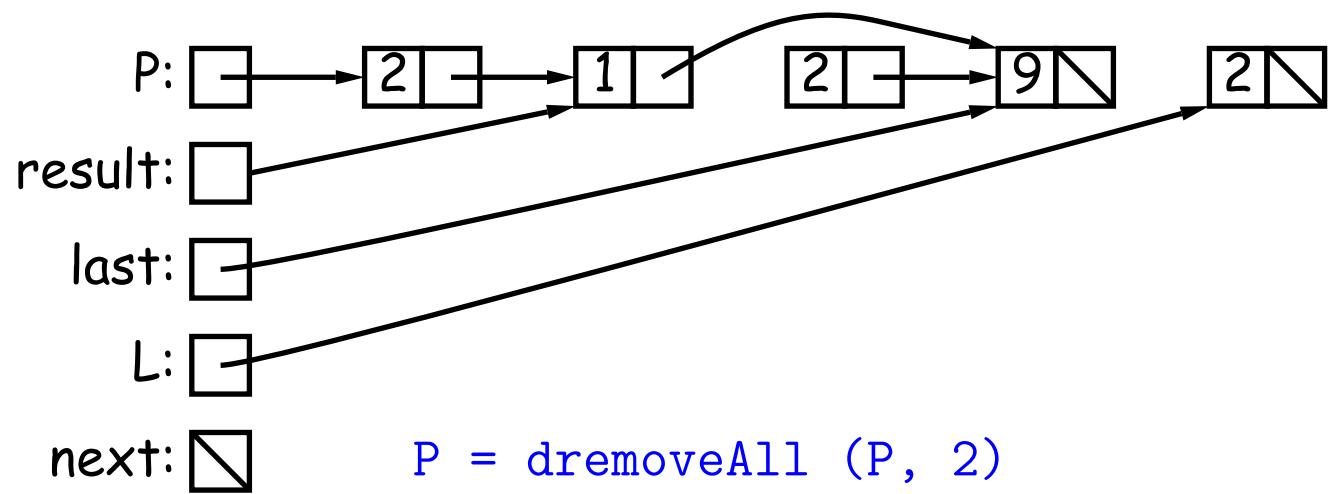
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



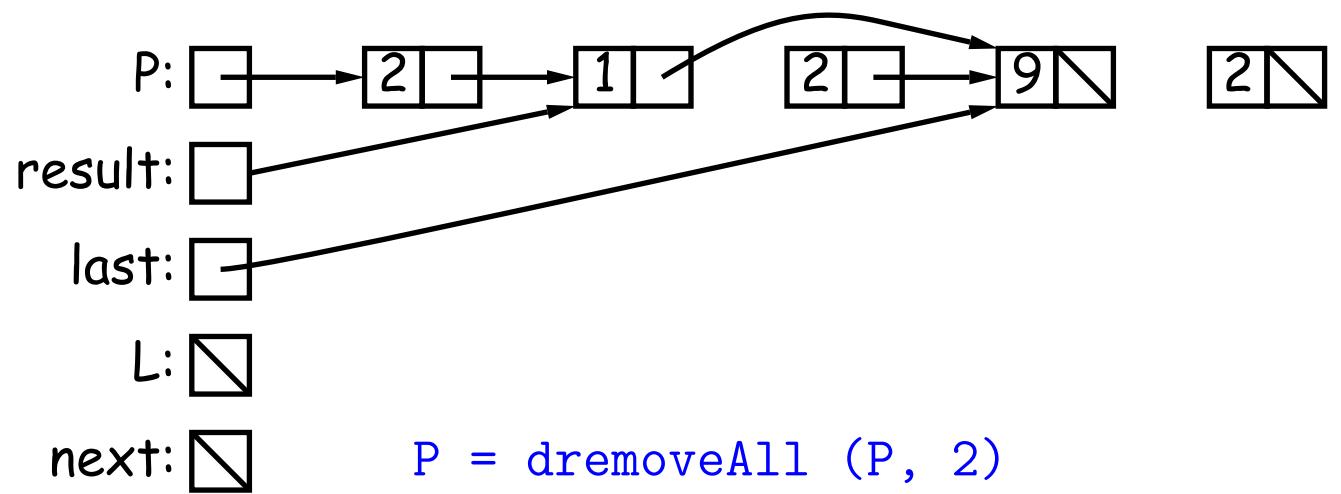
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



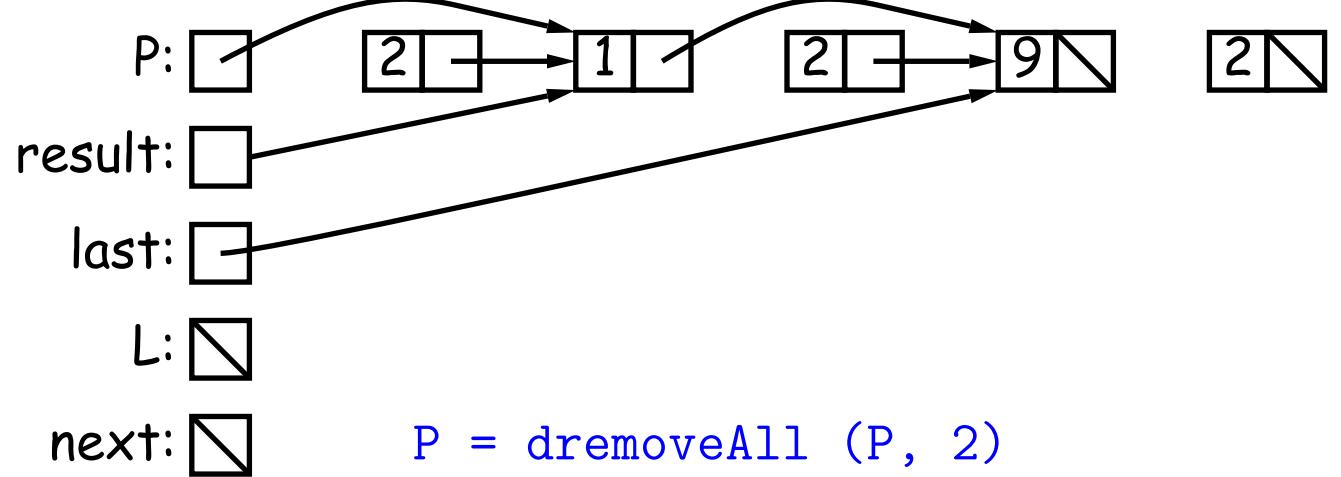
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



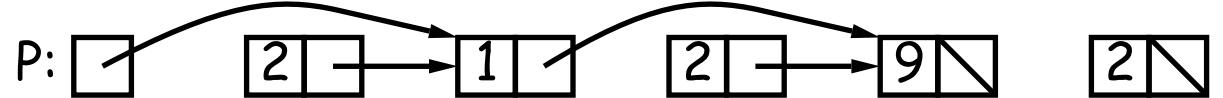
Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        // FIXME
    }
    return result;
}
```



P = dremoveAll (P, 2)

Iterative Destructive Deletion

```
/** The list resulting from removing all X's from L
 * destructively. */
static IntList dremoveAll(IntList L, int x) {
    IntList result, last, next;
    result = last = null;
    while (L != null) {
        next = L.tail;
        if (x != L.head) {
            if (last == null)
                result = last = L;
            else
                last = last.tail = L;
            L.tail = null;
        }
        L = next;
    }
    return result;
}
```



P = dremoveAll (P, 2)

Recreation

What is the sum of the coefficients of

$$(1 - 3x + 3x^2)^{743}(1 + 3x - 3x^2)^{744}$$

after expanding and collecting terms?

CS61B Lecture #5: Arrays

- An array is a structured container whose components are
 - `length`, a fixed integer.
 - a sequence of `length` simple containers of the same type, numbered from 0.
 - (`.length` field usually implicit in diagrams.)
- Arrays are anonymous, like other structured containers.
- Always referred to with pointers.
- For array pointed to by `A`,
 - Length is `A.length`
 - Numbered component i is `A[i]` (i is the `index`)
 - Important feature: index can be `any integer expression`.

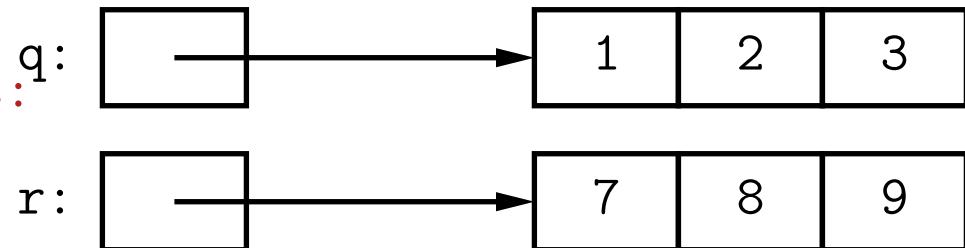
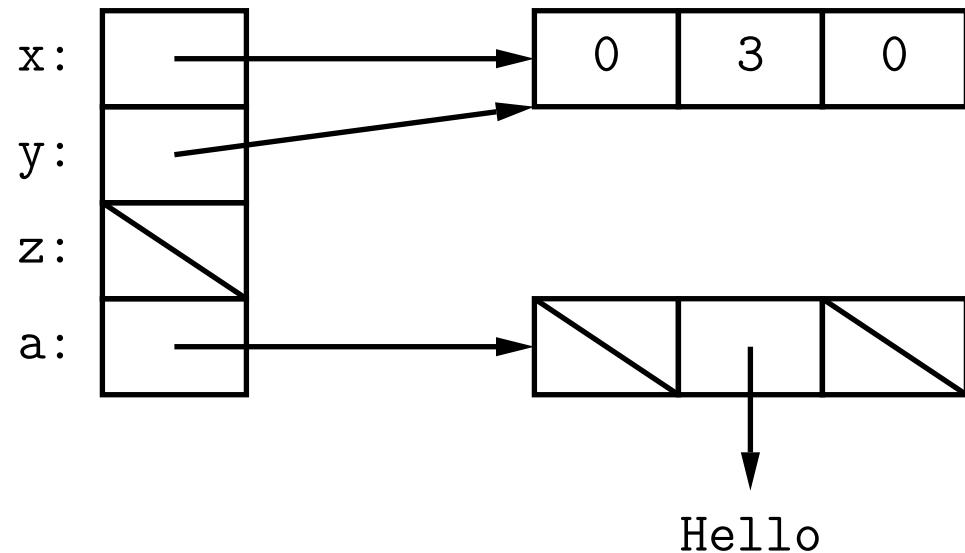
A Few Samples

Java

```
int[] x, y, z;  
String[] a;  
x = new int[3];  
y = x;  
a = new String[3];  
x[1] = 2;  
y[1] = 3;  
a[1] = "Hello";
```

```
int[] q;  
q = new int[] { 1, 2, 3 };  
// Short form for declarations:  
int[] r = { 7, 8, 9 };
```

Results



Example: Accumulate Values

Problem: Sum up the elements of array A.

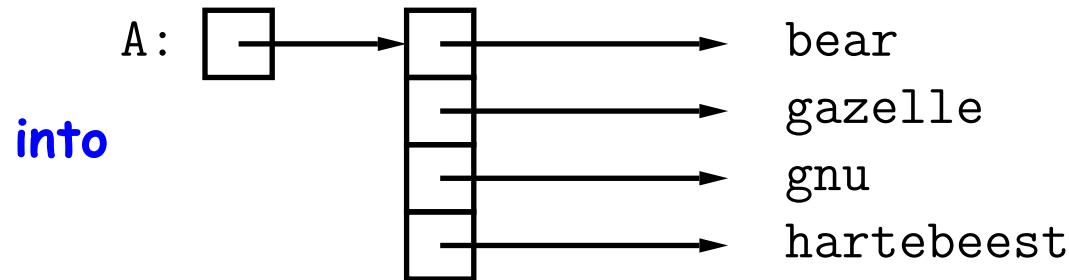
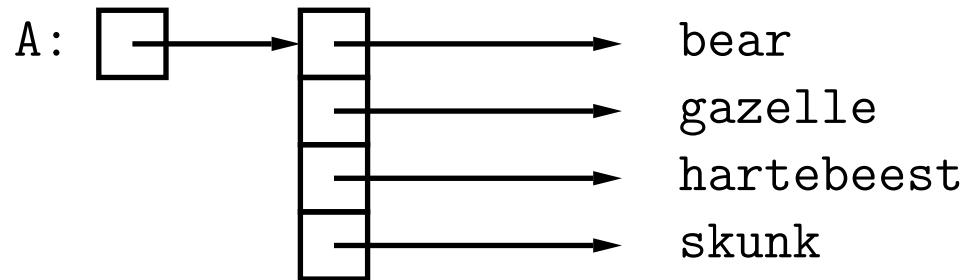
```
static int sum(int[] A) {  
    int N;  
    N = 0;                                // Newer syntax:  
    for (int i = 0; i < A.length; i += 1)  
        N += A[i];  
    return N;  
}  
  
// For the hard-core: could have written
```

```
int N, i;  
for (i=0, N=0; i<A.length; N += A[i], i += 1)  
    {} // or just ";" instead of "{}"
```

// But please don't: it's obscure.

Example: Insert into an Array

Problem: Want a call like `insert(A, 2, "gnu")` to convert (destructively)



into

```
/** Insert X at location K in ARR, moving items K, K+1, ... to locations
 * K+1, K+2, .... The last item in ARR is lost. */
static void insert(String[] arr, int k, String x) {
    for (int i = arr.length-1; i > k; i -= 1) // Why backwards?
        arr[i] = arr[i-1];
    /* Alternative to this loop:
       System.arraycopy(arr, k, arr, k+1, arr.length-k-1>);
               from      to      # to copy
    */
    arr[k] = x;
}
```

(Aside) Java Shortcut

- Useful tip: Can write just '`arraycopy`' by including at the top of the source file:

```
import static java.lang.System.arraycopy;
```

- This means "define the simple name `arraycopy` to be the equivalent of `java.lang.System.arraycopy` in the current source file."
- Can do the same for `out` so that you can write

```
out.println(...);
```

in place of

```
System.out.println(...);
```

- Finally, a declaration like

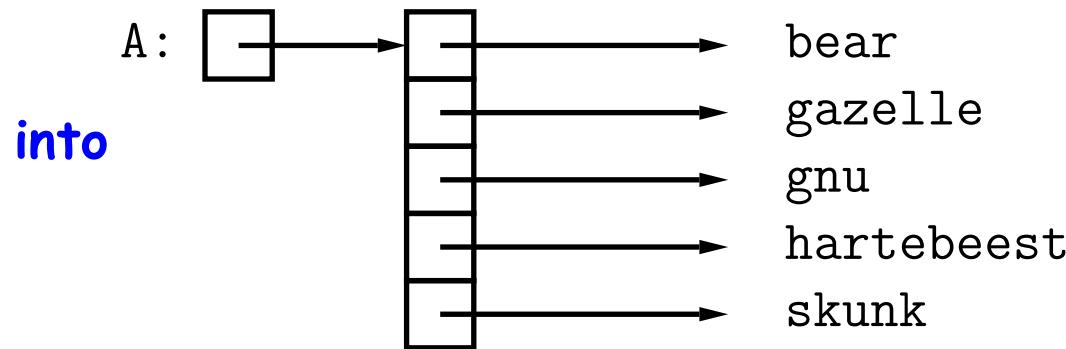
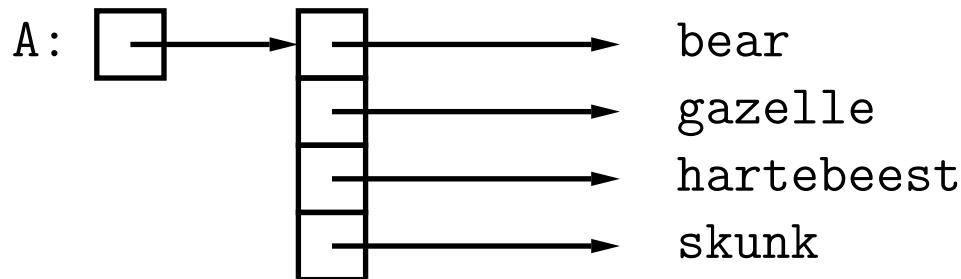
```
import static java.lang.Math.*;
```

means "take all the (public) static definitions in `java.lang.Math` and make them available in this source file by their simple names (the name after the last dot)."

- Useful for functions like `sin`, `sqrt`, etc.

Growing an Array

Problem: Suppose that we want to change the description above, so that $A = \text{insert2}(A, 2, "gnu")$ does *not* shove "skunk" off the end, but instead "grows" the array.

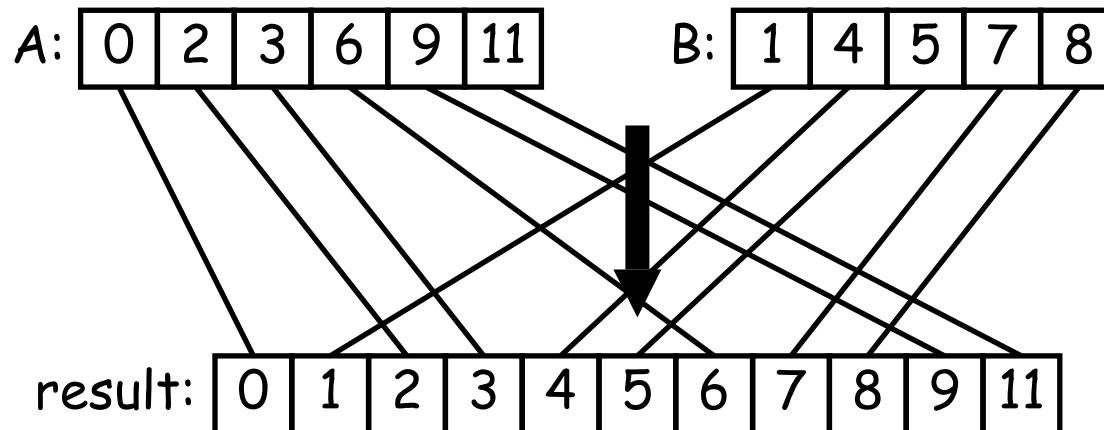


```
/** Return array, r, where r.length = ARR.length+1; r[0..K-1]
 *  the same as ARR[0..K-1], r[k] = x, r[K+1..] same as ARR[K..]. */
static String[] insert2(String[] arr, int k, String x) {
    String[] result = new String[arr.length + 1];
    arraycopy(arr, 0, result, 0, k);
    arraycopy(arr, k, result, k+1, arr.length-k);
    result[k] = x;
    return result;
}
```

Why do we need a different return type from insert2??

Example: Merging

Problem: Given two sorted arrays of ints, A and B, produce their **merge**: a sorted array containing all items from A and B.



Example: Merging Program

Problem: Given two sorted arrays of ints, A and B, produce their **merge**: a sorted array containing all from A and B.

Remark: In order to solve this recursively, it is useful to *generalize* the original function to allow merging *portions* of the arrays.

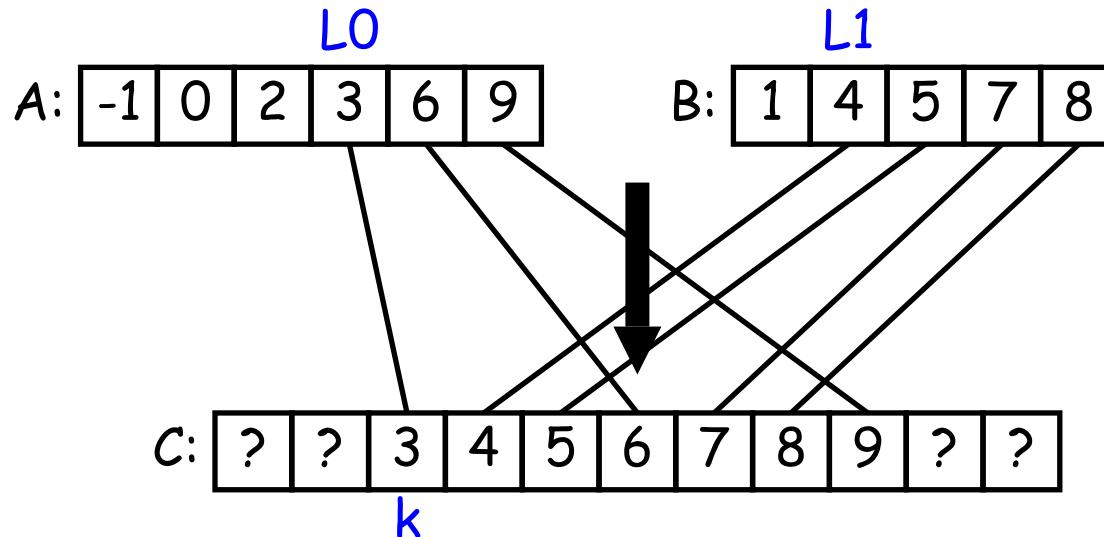
```
/** Assuming A and B are sorted, returns their merge. */
public static int[] merge(int[] A, int[] B) {
    return mergeTo(A, 0, B, 0);
}

/** The merge of A[L0..] and B[L1..] assuming A and B sorted. */
static int[] mergeTo(int[] A, int L0, int[] B, int L1) {
    int N = A.length - L0 + B.length - L1; int[] C = new int[N];
    if (L0 >= A.length) arraycopy(B, L1, C, 0, N);
    else if (L1 >= B.length) arraycopy(A, L0, C, 0, N);           What is wrong with
    else if (A[L0] <= B[L1]) {                                     this implementation?
        C[0] = A[L0]; arraycopy(mergeTo(A, L0+1, B, L1), 0, C, 1, N-1);
    } else {
        C[0] = B[L1]; arraycopy(mergeTo(A, L0, B, L1+1), 0, C, 1, N-1);
    }
    return C;
}
```

A Tail-Recursive Strategy

```
public static int[] merge(int[] A, int[] B) {  
    return mergeTo(A, 0, B, 0, new int[A.length+B.length], 0);  
}  
  
/** Merge A[L0..] and B[L1..] into C[K..], assuming A and B sorted. */  
static int[] mergeTo(int[] A, int L0, int[] B, int L1, int[] C, int k){  
    ...  
}
```

This last method merges *part* of A with part of B into part of C. For example, consider a possible call `mergeTo(A, 3, B, 1, C, 2)`



A Tail-Recursive Solution

```
public static int[] merge(int[] A, int[] B) {  
    return mergeTo(A, 0, B, 0, new int[A.length+B.length], 0);  
}  
  
/** Merge A[L0..] and B[L1..] into C[K..], assuming A and B sorted. */  
static int[] mergeTo(int[] A, int L0, int[] B, int L1, int[] C, int k){  
    if (??) {  
        return C;  
    } else if (??) {  
        C[k] = A[L0];  
        return mergeTo(A, ??, B, ??, C, ??)  
    } else {  
        C[k] = B[L1];  
        return mergeTo(A, ??, B, ??, C, ??)  
    }  
}
```

A Tail-Recursive Solution

```
public static int[] merge(int[] A, int[] B) {  
    return mergeTo(A, 0, B, 0, new int[A.length+B.length], 0);  
}  
  
/** Merge A[L0..] and B[L1..] into C[K..], assuming A and B sorted. */  
static int[] mergeTo(int[] A, int L0, int[] B, int L1, int[] C, int k){  
    if (L0 >= A.length && L1 >= B.length) {  
        return C;  
    } else if (??) {  
        C[k] = A[L0];  
        return mergeTo(A, ??, B, ??, C, ??)  
    } else {  
        C[k] = B[L1];  
        return mergeTo(A, ??, B, ??, C, ??)  
    }  
}
```

A Tail-Recursive Solution

```
public static int[] merge(int[] A, int[] B) {  
    return mergeTo(A, 0, B, 0, new int[A.length+B.length], 0);  
}  
  
/** Merge A[L0..] and B[L1..] into C[K..], assuming A and B sorted. */  
static int[] mergeTo(int[] A, int L0, int[] B, int L1, int[] C, int k){  
    if (L0 >= A.length && L1 >= B.length) {  
        return C;  
    } else if (L1 >= B.length || (L0 < A.length && A[L0] <= B[L1])) {  
        C[k] = A[L0];  
        return mergeTo(A, ??, B, ??, C, ??)  
    } else {  
        C[k] = B[L1];  
        return mergeTo(A, ??, B, ??, C, ??)  
    }  
}
```

A Tail-Recursive Solution

```
public static int[] merge(int[] A, int[] B) {  
    return mergeTo(A, 0, B, 0, new int[A.length+B.length], 0);  
}  
  
/** Merge A[L0..] and B[L1..] into C[K..], assuming A and B sorted. */  
static int[] mergeTo(int[] A, int L0, int[] B, int L1, int[] C, int k){  
    if (L0 >= A.length && L1 >= B.length) {  
        return C;  
    } else if (L1 >= B.length || (L0 < A.length && A[L0] <= B[L1])) {  
        C[k] = A[L0];  
        return mergeTo(A, L0 + 1, B, L1, C, k + 1);  
    } else {  
        C[k] = B[L1];  
        return mergeTo(A, ??, B, ??, C, ??)  
    }  
}
```

A Tail-Recursive Solution

```
public static int[] merge(int[] A, int[] B) {  
    return mergeTo(A, 0, B, 0, new int[A.length+B.length], 0);  
}  
  
/** Merge A[L0..] and B[L1..] into C[K..], assuming A and B sorted. */  
static int[] mergeTo(int[] A, int L0, int[] B, int L1, int[] C, int k){  
    if (L0 >= A.length && L1 >= B.length) {  
        return C;  
    } else if (L1 >= B.length || (L0 < A.length && A[L0] <= B[L1])) {  
        C[k] = A[L0];  
        return mergeTo(A, L0 + 1, B, L1, C, k + 1);  
    } else {  
        C[k] = B[L1];  
        return mergeTo(A, L0, B, L1 + 1, C, k + 1);  
    }  
}
```

Iterative Solution

In general, we don't use either of the previous approaches in languages like C and Java. Array manipulation is most often iterative:

```
public static int[] merge(int[] A, int[] B) {  
    int[] C = new int[A.length + B.length];  
    // mergeTo(A, 0, B, 0, C, 0)  
    int L0, L1, k;  
    L0 = L1 = k = 0;  
  
    while (??) {  
        if (L1 >= B.length || (L0 < A.length && A[L0] <= B[L1])) {  
            C[k] = A[L0];  
            ??  
        } else {  
            C[k] = B[L1];  
            ??  
        }  
    }  
    return C;  
}
```

Iterative Solution

In general, we don't use either of the previous approaches in languages like C and Java. Array manipulation is most often iterative:

```
public static int[] merge(int[] A, int[] B) {
    int[] C = new int[A.length + B.length];
    // mergeTo(A, 0, B, 0, C, 0)
    int L0, L1, k;
    L0 = L1 = k = 0;

    while (L0 < A.length || L1 < B.length) {
        if (L1 >= B.length || (L0 < A.length && A[L0] <= B[L1])) {
            C[k] = A[L0];
            ???
        } else {
            C[k] = B[L1];
            ???
        }
    }
    return C;
}
```

Iterative Solution

In general, we don't use either of the previous approaches in languages like C and Java. Array manipulation is most often iterative:

```
public static int[] merge(int[] A, int[] B) {  
    int[] C = new int[A.length + B.length];  
    // mergeTo(A, 0, B, 0, C, 0)  
    int L0, L1, k;  
    L0 = L1 = k = 0;  
  
    while (L0 < A.length || L1 < B.length) {  
        if (L1 >= B.length || (L0 < A.length && A[L0] <= B[L1])) {  
            C[k] = A[L0];  
            L0 += 1; k += 1;  
        } else {  
            C[k] = B[L1];  
            L1 += 1; k += 1;  
        }  
    }  
    return C;  
}
```

Iterative Solution II

The same, with a `for` loop:

```
public static int[] merge(int[] A, int[] B) {  
    int[] C = new int[A.length + B.length];  
    int L0, L1;  
    L0 = L1 = 0;  
    for (int k = 0; k < C.length; k += 1) {  
        if (L1 >= B.length || (L0 < A.length && A[L0] <= B[L1])) {  
            C[k] = A[L0]; L0 += 1;  
        } else {  
            C[k] = B[L1]; L1 += 1;  
        }  
    }  
    return C;  
}
```

Observation: k is redundant, because $k = ?$.

Iterative Solution II

The same, with a `for` loop:

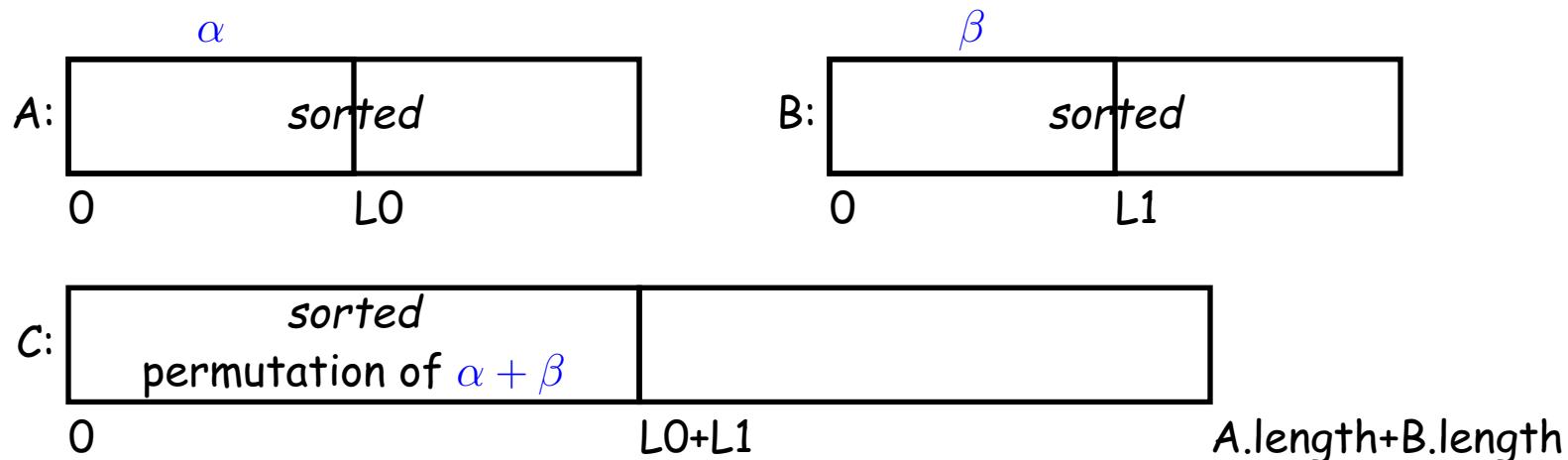
```
public static int[] merge(int[] A, int[] B) {  
    int[] C = new int[A.length + B.length];  
    int L0, L1;  
    L0 = L1 = 0;  
    for (int k = 0; k < C.length; k += 1) {  
        if (L1 >= B.length || (L0 < A.length && A[L0] <= B[L1])) {  
            C[k] = A[L0]; L0 += 1;  
        } else {  
            C[k] = B[L1]; L1 += 1;  
        }  
    }  
    return C;  
}
```

Observation: k is redundant, because $k = L0 + L1$.

Alternative Solution: Removing k

Using previous invariant that $k=L_0+L_1$ simplifies things:

```
public static int[] merge(int[] A, int[] B) {  
    int[] C = new int[A.length + B.length];  
    int L0, L1; L0 = L1 = 0;  
    while (L0 + L1 < C.length) {  
        if (L1 >= B.length || (L0 < A.length && A[L0] < B[L1])) {  
            C[L0 + L1] = A[L0]; L0 += 1;  
        } else {  
            C[L0 + L1] = B[L1]; L1 += 1;  
        }  
    }  
    return C;  
}
```



Multidimensional Arrays

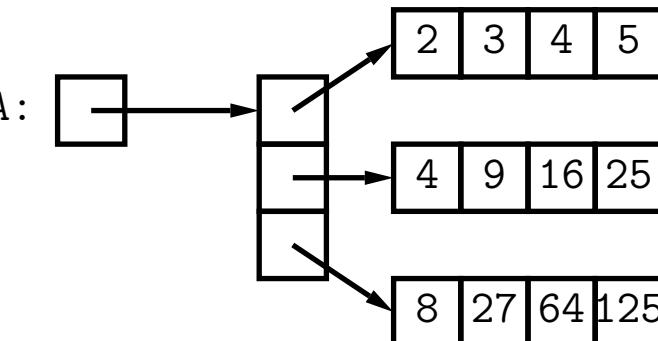
What about two- or higher-dimensional layouts, such as

$$A = \begin{array}{|c|c|c|c|} \hline 2 & 3 & 4 & 5 \\ \hline 4 & 9 & 16 & 25 \\ \hline 8 & 27 & 64 & 125 \\ \hline \end{array} \quad ?$$

Multidimensional Arrays in Java

These are not primitive in Java, but we can build them as [arrays of arrays](#):

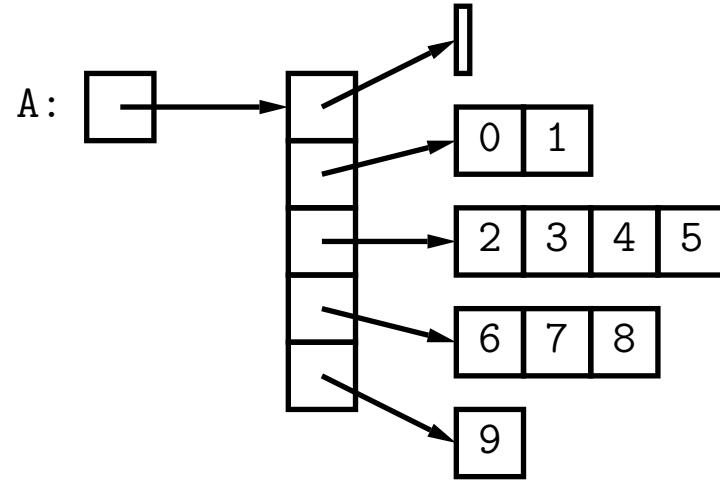
```
int[][] A = new int[3][];  
A[0] = new int[] {2, 3, 4, 5};  
A[1] = new int[] {4, 9, 16, 25};  
A[2] = new int[] {8, 27, 64, 125};  
  
// or  
  
int[][] A;  
A = new int[][] { {2, 3, 4, 5},  
                 {4, 9, 16, 25},  
                 {8, 27, 64, 125} };  
  
// or  
  
int[][] A = { {2, 3, 4, 5},  
              {4, 9, 16, 25},  
              {8, 27, 64, 125} };  
  
// or  
  
int[][] A = new int[3][4];  
for (int i = 0; i < 3; i += 1)  
    for (int j = 0; j < 4; j += 1)  
        A[i][j] = (int) Math.pow(j + 2, i + 1);
```



Exotic Multidimensional Arrays

- Since every element of an array is independent, there is no single "width" in general:

```
int [] [] A = new int [5] [] ;  
A[0] = new int [] {} ;  
A[1] = new int [] {0, 1} ;  
A[2] = new int [] {2, 3, 4, 5} ;  
A[3] = new int [] {6, 7, 8} ;  
A[4] = new int [] {9} ;
```



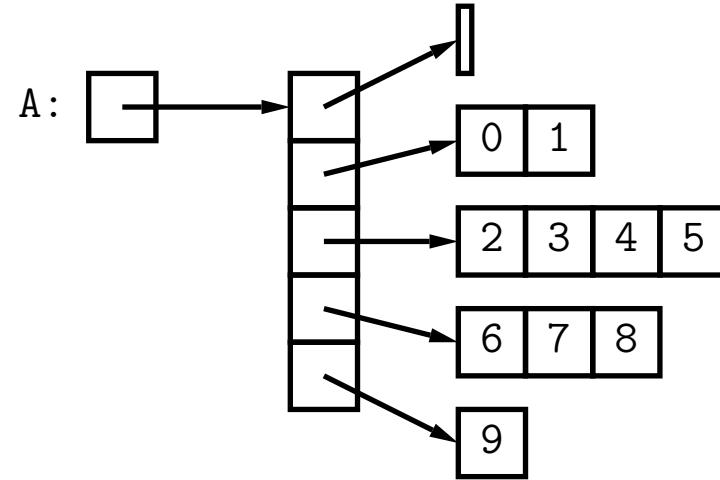
- What does this print?

```
int [] [] ZERO = new int [3] [] ;  
ZERO[0] = ZERO[1] = ZERO[2] =  
    new int [] {0, 0, 0} ;  
ZERO[0][1] = 1 ;  
System.out.println(ZERO[2][1]) ;
```

Exotic Multidimensional Arrays

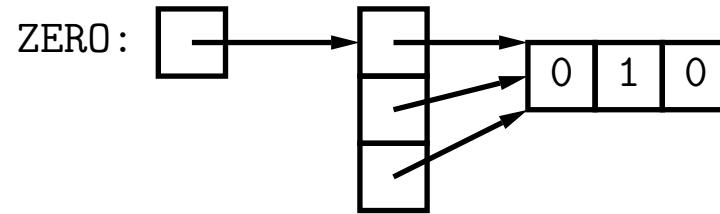
- Since every element of an array is independent, there is no single "width" in general:

```
int [] [] A = new int [5] [] ;  
A[0] = new int [] {} ;  
A[1] = new int [] {0, 1} ;  
A[2] = new int [] {2, 3, 4, 5} ;  
A[3] = new int [] {6, 7, 8} ;  
A[4] = new int [] {9} ;
```



- What does this print?

```
int [] [] ZERO = new int [3] [] ;  
ZERO[0] = ZERO[1] = ZERO[2] =  
    new int [] {0, 0, 0} ;  
ZERO[0][1] = 1 ;  
System.out.println(ZERO[2][1]) ;
```



CS61B Lecture #6: More Iteration: Sort an Array

Problem. Print out the command-line arguments in lexicographic order:

```
% java Sort the quick brown fox jumped over the lazy dog  
brown dog fox jumped lazy over quick the the
```

Plan.

```
public class Sort {  
    /** Sort and print WORDS lexicographically. */  
    public static void main(String[] words) {  
        sort(words, 0, words.length-1);  
        print(words);  
    }  
  
    /** Sort items A[L..U] , with all others unchanged. */  
    static void sort(String[] A, int L, int U) { /* "TOMORROW" */ }  
  
    /** Print A on one line, separated by blanks. */  
    static void print(String[] A) { /* "TOMORROW" */ }  
}
```

How do We Know If It Works?

- *Unit testing* refers to the testing of individual units (methods) within a program, rather than the whole program.
- *Module testing* refers to testing of classes or other groupings of methods and data.
- *System testing* (or *acceptance testing*) refers to the testing of the functionality of an entire program.
- *Integration testing* is sort of intermediate between unit and system testing, and tests that modules work correctly together.
- *Regression testing* refers to testing with the specific goal of checking that fixes, enhancements, or other changes have not introduced faults (regressions).
- In this class, we mainly use the JUnit tool for unit testing.
- Example: AGTestYear.java in lab #1.
- Our system testing is somewhat more *ad hoc*, and customized to the program. At its simplest, one might just run specific input files through the program and compare with precomputed outputs.

Test-Driven Development

- Idea: write tests first.
- Implement unit at a time, run tests, fix and refactor until it works.
- We're only going to push is fairly lightly in this course, but it is useful and has quite a following.
- You'll see a lot more of it in CS169.

Testing sort

- This is pretty easy: just give a bunch of arrays to sort and then make sure they each get sorted properly.
- Have to make sure we cover the necessary cases:
 - *Corner cases.* E.g., empty array, one-element, all elements the same.
 - *Representative "middle" cases.* E.g., elements reversed, elements in order, one pair of elements reversed,

Simple JUnit

- The JUnit package provides some handy tools for unit testing.
- The Java annotation `@Test` on a method tells the JUnit machinery to call that method.
- (An *annotation* in Java provides information about a method, class, etc., that can be examined within Java itself.)
- A collection of methods with names beginning with `assert` then allow your test cases to check conditions and report failures.
- [See example in the [code](#) link for lecture 6.]

Selection Sort

```
/** Sort items A[L..U], with all others unchanged. */
static void sort(String[] A, int L, int U) {
    if (L < U) {
        int k = /*( Index s.t. A[k] is largest in A[L], ..., A[U] )*/;
        /*{ swap A[k] with A[U] }*/;
        /*{ Sort items L to U-1 of A. }*/;
    }
}
```

And we're done! Well, OK, not quite.

Selection Sort

```
/** Sort items A[L..U], with all others unchanged. */
static void sort(String[] A, int L, int U) {
    if (L < U) {
        int k = indexOfLargest(A, L, U);
        /*{ swap A[k] with A[U] }*/;
        /*{ Sort items L to U-1 of A. }*/;
    }
}

/** Index k, I0<=k<=I1, such that V[k] is largest element among
 * V[I0], ... V[I1]. Requires I0<=I1. */
static int indexOfLargest(String[] V, int i0, int i1) {
    ...
}
```

Selection Sort

```
/** Sort items A[L..U], with all others unchanged. */
static void sort(String[] A, int L, int U) {
    if (L < U) {
        int k = indexOfLargest(A, L, U);
        /*{ swap A[k] with A[U] }*/;
        sort(A, L, U-1);           // Sort items L to U-1 of A
    }
}

/** Index k, I0<=k<=I1, such that V[k] is largest element among
 * V[I0], ... V[I1]. Requires I0<=I1. */
static int indexOfLargest(String[] V, int i0, int i1) {
    ...
}
```

Selection Sort

```
/** Sort items A[L..U], with all others unchanged. */
static void sort(String[] A, int L, int U) {
    if (L < U) {
        int k = indexOfLargest(A, L, U);
        String tmp = A[k]; A[k] = A[U]; A[U] = tmp;
        sort(A, L, U-1);           // Sort items L to U-1 of A
    }
}

/** Index k, I0<=k<=I1, such that V[k] is largest element among
 * V[I0], ... V[I1]. Requires I0<=I1. */
static int indexOfLargest(String[] V, int i0, int i1) {
    ...
}
```

Selection Sort

```
/** Sort items A[L..U], with all others unchanged. */
static void sort(String[] A, int L, int U) {
    if (L < U) {
        int k = indexOfLargest(A, L, U);
        String tmp = A[k]; A[k] = A[U]; A[U] = tmp;
        sort(A, L, U-1);           // Sort items L to U-1 of A
    }
}
```

What would an iterative version look like?

```
while (?) {
    ?
}
```

Selection Sort

```
/** Sort items A[L..U], with all others unchanged. */
static void sort(String[] A, int L, int U) {
    if (L < U) {
        int k = indexOfLargest(A, L, U);
        String tmp = A[k]; A[k] = A[U]; A[U] = tmp;
        sort(A, L, U-1);           // Sort items L to U-1 of A
    }
}
```

Iterative version:

```
while (L < U) {
    int k = indexOfLargest(A, L, U);
    String tmp = A[k]; A[k] = A[U]; A[U] = tmp;
    U -= 1;
}
```

Find Largest

```
/** Index k, I0<=k<=I1, such that V[k] is largest element among
 * V[I0], ... V[I1]. Requires I0<=I1. */
static int indexOfLargest(String[] V, int i0, int i1) {
    if (?)
        return i1;
    else {
        }
}
```

Find Largest

```
/** Index k, I0<=k<=I1, such that V[k] is largest element among
 * V[I0], ... V[I1]. Requires I0<=I1. */
static int indexOfLargest(String[] V, int i0, int i1) {
    if (i0 >= i1)
        return i1;
    else /* if (i0 < i1) */ {
        }
}
```

Find Largest

```
/** Index k, I0<=k<=I1, such that V[k] is largest element among
 * V[I0], ... V[I1]. Requires I0<=I1. */
static int indexOfLargest(String[] V, int i0, int i1) {
    if (i0 >= i1)
        return i1;
    else /* if (i0 < i1) */ {
        int k = /*( index of largest value in V[i0 + 1..i1] );*/
        return /*( whichever of i0 and k has larger value );*/
    }
}
```

Find Largest

```
/** Index k, I0<=k<=I1, such that V[k] is largest element among
 * V[I0], ... V[I1]. Requires I0<=I1. */
static int indexOfLargest(String[] V, int i0, int i1) {
    if (i0 >= i1)
        return i1;
    else /* if (i0 < i1) */ {
        int k = indexOfLargest(V, i0 + 1, i1);
        return /*( whichever of i0 and k has larger value )*/;
    }
}
```

Find Largest

```
/** Index k, I0<=k<=I1, such that V[k] is largest element among
 * V[I0], ... V[I1]. Requires I0<=I1. */
static int indexOfLargest(String[] V, int i0, int i1) {
    if (i0 >= i1)
        return i1;
    else /* if (i0 < i1) */ {
        int k = indexOfLargest(V, i0 + 1, i1);
        return (V[i0].compareTo(V[k]) > 0) ? i0 : k;
        // if (v[i0].compareTo(v[k]) > 0) return i0; else return k;
    }
}
```

- Turning this into an iterative version is tricky: not tail recursive.
- What are the arguments to compareTo the first time it's called?

Iteratively Find Largest

```
/** Value k, I0<=k<=I1, such that V[k] is largest element among
 * V[I0], ... V[I1]. Requires I0<=I1. */
static int indexOfLargest(String[] V, int i0, int i1) {
    if (i0 >= i1)
        return i1;
    else /* if (i0 < i1) */ {
        int k = indexOfLargest(V, i0 + 1, i1);
        return (V[i0].compareTo(V[k]) > 0) ? i0 : k;
        // if (V[i0].compareTo(V[k]) > 0) return i0; else return k;
    }
}
```

Iterative:

```
int i, k;
k = ?;      // Deepest iteration
for (i = ?; ...?; i ...?)
    k = ?;
return k;
```

Iteratively Find Largest

```
/** Value k, I0<=k<=I1, such that V[k] is largest element among
 * V[I0], ... V[I1]. Requires I0<=I1. */
static int indexOfLargest(String[] V, int i0, int i1) {
    if (i0 >= i1)
        return i1;
    else /* if (i0 < i1) */ {
        int k = indexOfLargest(V, i0 + 1, i1);
        return (V[i0].compareTo(V[k]) > 0) ? i0 : k;
        // if (V[i0].compareTo(V[k]) > 0) return i0; else return k;
    }
}
```

Iterative:

```
int i, k;
k = i1];      // Deepest iteration
for (i = ?; ...?; i ...?)
    k = ?;
return k;
```

Iteratively Find Largest

```
/** Value k, I0<=k<=I1, such that V[k] is largest element among
 * V[I0], ... V[I1]. Requires I0<=I1. */
static int indexOfLargest(String[] V, int i0, int i1) {
    if (i0 >= i1)
        return i1;
    else /* if (i0 < i1) */ {
        int k = indexOfLargest(V, i0 + 1, i1);
        return (V[i0].compareTo(V[k]) > 0) ? i0 : k;
        // if (V[i0].compareTo(V[k]) > 0) return i0; else return k;
    }
}
```

Iterative:

```
int i, k;
k = i1];      // Deepest iteration
for (i = i1 - 1; i >= i0; i -= 1)
    k = ?;
return k;
```

Iteratively Find Largest

```
/** Value k, I0<=k<=I1, such that V[k] is largest element among
 * V[I0], ... V[I1]. Requires I0<=I1. */
static int indexOfLargest(String[] V, int i0, int i1) {
    if (i0 >= i1)
        return i1;
    else /* if (i0 < i1) */ {
        int k = indexOfLargest(V, i0 + 1, i1);
        return (V[i0].compareTo(V[k]) > 0) ? i0 : k;
        // if (V[i0].compareTo(V[k]) > 0) return i0; else return k;
    }
}
```

Iterative:

```
int i, k;
k = i1;      // Deepest iteration
for (i = i1 - 1; i >= i0; i -= 1)
    k = (V[i].compareTo(V[k]) > 0) ? i : k;
return k;
```

Finally, Printing

```
/** Print A on one line, separated by blanks. */
static void print(String[] A) {
    for (int i = 0; i < A.length; i += 1)
        System.out.print(A[i] + " ");
    System.out.println();
}

/* Java also provides a simple, specialized syntax for looping
 * through an entire array: */
for (String s : A)
    System.out.print(s + " ");
```

Another Problem

Given an array of integers, A , of length $N > 0$, find the smallest index, k , such that all elements at indices $\geq k$ and $< N - 1$ are greater than $A[N - 1]$. Then rotate elements k to $N - 1$ right by one. For example, if A starts out as

$\{ 1, 9, 4, 3, 0, 12, 11, 9, 15, 22, 12 \}$

then it ends up as

$\{ 1, 9, 4, 3, 0, 12, 11, 9, 12, 15, 22 \}$

As another example,

$\{ 1, 9, 4, 3, 0, 12, 11, 9, 15, 22, -2 \}$

would become

$\{ -2, 1, 9, 4, 3, 0, 12, 11, 9, 15, 22 \}$

What if A starts like this?

$\{ 1, 9, 4, 3, 0, 12, 11, 9, 12, 15, 22 \}$

Another Problem

Given an array of integers, A , of length $N > 0$, find the smallest index, k , such that all elements at indices $\geq k$ and $< N - 1$ are greater than $A[N - 1]$. Then rotate elements k to $N - 1$ right by one. For example, if A starts out as

$\{ 1, 9, 4, 3, 0, 12, 11, 9, 15, 22, 12 \}$

then it ends up as

$\{ 1, 9, 4, 3, 0, 12, 11, 9, 12, 15, 22 \}$

As another example,

$\{ 1, 9, 4, 3, 0, 12, 11, 9, 15, 22, -2 \}$

would become

$\{ -2, 1, 9, 4, 3, 0, 12, 11, 9, 15, 22 \}$

What if A starts like this?

$\{ 1, 9, 4, 3, 0, 12, 11, 9, 12, 15, 22 \}$

Answer: It's unchanged. (No, the spec is not ambiguous.)

Your turn

```
public class Shove {  
  
    /** Rotate elements A[k] to A[A.length-1] one element to the  
     * right, where k is the smallest index such that elements  
     * k through A.length-2 are all larger than A[A.length-1].  
     */  
    static void moveOver(int[] A) {  
        // FILL IN  
    }  
}
```

Recreation

Given that

$$\log(1+x) = x - \frac{1}{2}x^2 + \frac{1}{3}x^3 - \dots$$

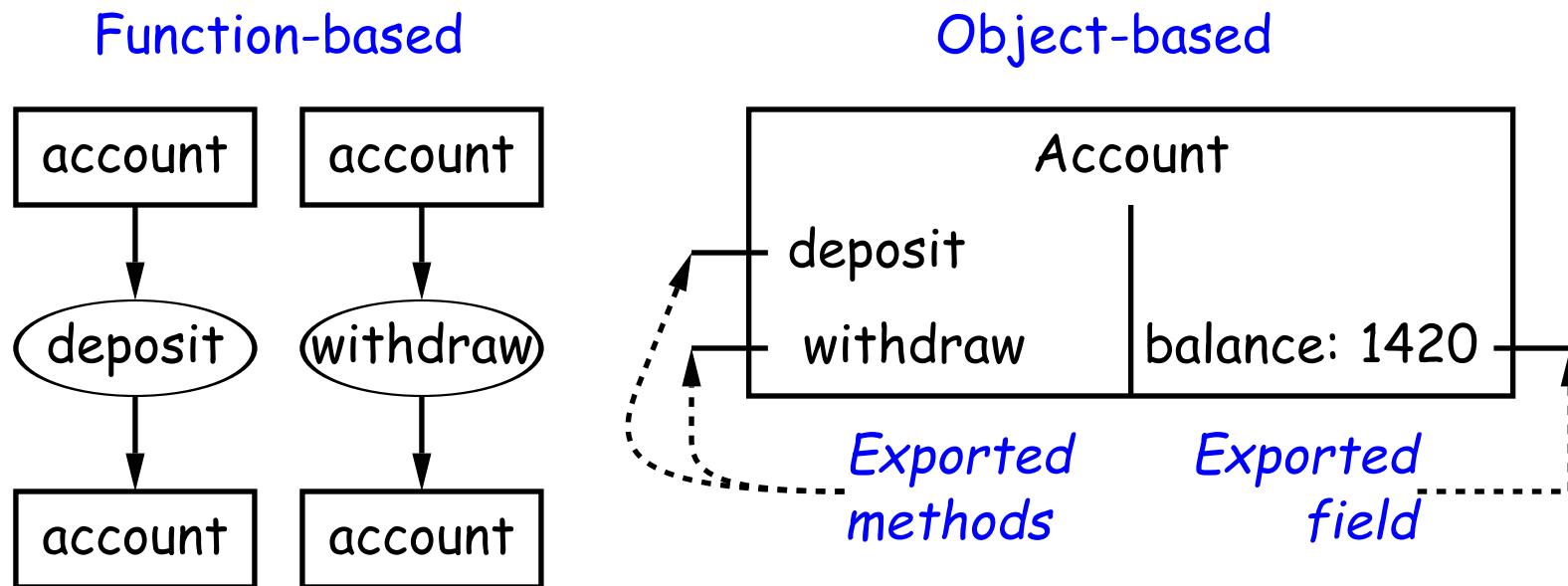
why is it not the case that

$$\begin{aligned}\log 2 &= 1 - 1/2 + 1/3 - 1/4 + 1/5 - 1/6 + 1/7 - 1/8 + 1/9 - \dots \\&= (1 + 1/3 + 1/5 + 1/7 + 1/9 + \dots) - (1/2 + 1/4 + 1/6 + 1/8 + \dots) \\&= (1 + 1/3 + 1/5 + 1/7 + 1/9 + \dots) + (1/2 + 1/4 + 1/6 + 1/8 + \dots) \\&\quad - 2(1/2 + 1/4 + 1/6 + 1/8 + \dots) \\&= (1 + 1/2 + 1/3 + 1/4 + \dots) - (1 + 1/2 + 1/3 + 1/4 + \dots) \\&= 0?\end{aligned}$$

CS61B Lecture #7: Object-Based Programming

Basic Idea.

- *Function-based programs* are organized primarily around the functions (methods, etc.) that do things. Data structures (objects) are considered separate.
- *Object-based programs* are organized around the *types of objects* that are used to represent data; methods are grouped by type of object.
- Simple banking-system example:



Philosophy

- Idea (from 1970s and before): An *abstract data type* is
 - a set of possible values (a *domain*), plus
 - a set of *operations* on those values (or their containers).
- In IntList, for example, the domain was a *set of pairs*: (head, tail), where head is an int and tail is a pointer to an IntList.
- The IntList operations consisted only of assigning to and accessing the two fields (head and tail).
- In general, we prefer a purely *procedural interface*, where the functions (methods) do everything—no outside access to the internal representation (i.e., instance variables).
- That way, implementor of a class and its methods has complete control over behavior of instances.
- In Java, the preferred way to write the “operations of a type” is as *instance methods*.

You Saw It All (Maybe) in CS61A: The Account Class

```
class Account:  
    balance = 0  
    def __init__(self, balance0):  
        self.balance = balance0  
  
    def deposit(self, amount):  
        self.balance += amount  
        return self.balance  
  
    def withdraw(self, amount):  
        if self.balance < amount:  
            raise ValueError \  
                ("Insufficient funds")  
        else:  
            self.balance -= amount  
        return self.balance
```

```
myAccount = Account(1000)  
print(myAccount.balance)  
myAccount.deposit(100)  
myAccount.withdraw(500)
```

```
public class Account {  
    public int balance;  
    public Account(int balance0) {  
        this.balance = balance0;  
    }  
    public int deposit(int amount) {  
        balance += amount; return balance;  
    }  
    public int withdraw(int amount) {  
        if (balance < amount)  
            throw new IllegalStateException  
                ("Insufficient funds");  
        else balance -= amount;  
        return balance;  
    }  
}
```

```
Account myAccount = new Account(1000);  
print(myAccount.balance)  
myAccount.deposit(100);  
myAccount.withdraw(500);
```

You Also Would Have Seen It All in CS61AS

```
(define-class (account balance0)
  (instance-vars (balance 0))
  (initialize
    (set! balance balance0))

  (method (deposit amount)
    (set! balance (+ balance amount))
    balance)
  (method (withdraw amount)
    (if (< balance amount)
        (error "Insufficient funds")
        (begin
          (set! balance (- balance amount))
          balance)))))



---


(define my-account
  (instantiate account 1000))
(ask my-account 'balance)
(ask my-account 'deposit 100)
(ask my-account 'withdraw 500)
```

```
public class Account {
  public int balance;
  public Account(int balance0) {
    balance = balance0;
  }
  public int deposit(int amount) {
    balance += amount; return balance;
  }
  public int withdraw(int amount) {
    if (balance < amount)
      throw new IllegalStateException
        ("Insufficient funds");
    else balance -= amount;
    return balance;
  }
}
```

```
Account myAccount = new Account(1000);
myAccount.balance
myAccount.deposit(100);
myAccount.withdraw(500);
```

The Pieces

- **Class declaration** defines a *new type of object*, i.e., new type of structured container.
- **Instance variables** such as balance are the simple containers within these objects (*fields* or *components*).
- **Instance methods**, such as deposit and withdraw are like ordinary (static) methods that take an invisible extra parameter (called **this**).
- The **new** operator creates (*instantiates*) new objects, and initializes them using constructors.
- **Constructors** such as the method-like declaration of Account are special methods that are used only to initialize new instances. They take their arguments from the **new** expression.
- **Method selection** picks methods to call. For example,

```
myAccount.deposit(100)
```

tells us to call the method named deposit that is defined for the object pointed to by myAccount.

Getter Methods

- Slight problem with Java version of Account: anyone can assign to the balance field
- This reduces the control that the implementor of Account has over possible values of the balance.
- Solution: allow public access only through methods:

```
public class Account {  
    private int _balance;  
  
    ...  
    public int balance() { return _balance; }  
    ...  
}
```

- Now Account._balance = 1000000 is an error outside Account.
- (I use the convention of putting '_' at the start of private instance variables to distinguish them from local variables and non-private variables. Could actually use balance for both the method and the variable, but please don't.)

Class Variables and Methods

- Suppose we want to keep track of the bank's total funds.
- This number is not associated with any particular Account, but is common to all—it is *class-wide*. In Java, “class-wide” \equiv static.

```
public class Account {  
    ...  
    private static int _funds = 0;  
    public int deposit(int amount) {  
        _balance += amount;  
        _funds += amount;          // or this._funds or Account._funds  
        return _balance;  
    }  
    public static int funds() {  
        return _funds;            // or Account._funds  
    }  
    ... // Also change withdraw.  
}
```

- From outside, can refer to either Account.funds() or to myAccount.funds() (same thing).

Instance Methods

- Instance method such as

```
int deposit(int amount) {  
    _balance += amount;  
    _funds += amount;  
    return balance;  
}
```

behaves sort of like a static method with hidden argument:

```
static int deposit(final Account this, int amount) {  
    this._balance += amount;  
    _funds += amount;  
    return this._balance;  
}
```

- NOTE: Just explanatory: Not real Java (not allowed to declare 'this'). (final *is* real Java; means "can't change once initialized.")

Calling Instance Methods

```
/** (Fictional) equivalent of deposit instance method. */
static int deposit(final Account this, int amount) {
    this._balance += amount;
    _funds += amount;
    return this._balance;
}
```

- Likewise, the instance-method call `myAccount.deposit(100)` is like a call on this fictional static method:

```
Account.deposit(myAccount, 100);
```

- Compare `this` with Python, where the extra parameter is not fictional (and one can choose the name—usually `self`.)
- Inside a real instance method, as a convenient abbreviation, one can leave off the leading '`this.`' on field access or method call if not ambiguous. (Unlike Python)

'Instance' and 'Static' Don't Mix

- Since real static methods don't have the invisible this parameter, makes no sense to refer directly to instance variables in them:

```
public static int badBalance(Account A) {  
    int x = A._balance; // This is OK  
                           // (A tells us whose balance)  
  
    return _balance; // WRONG! NONSENSE!  
}
```

- Reference to _balance here equivalent to this._balance,
- But this is meaningless (whose balance?)
- However, it makes perfect sense to access a static (class-wide) field or method in an instance method or constructor, as happened with _funds in the deposit method.
- There's only one of each static field, so don't need to have a 'this' to get it. Can just name the class (or use no qualification inside the class, as we've been doing).

Constructors

- To completely control objects of some class, you must be able to set their initial contents.
- A *constructor* is a kind of special instance method that is called by the **new** operator right after it creates a new object, as if

`L = new IntList(1,null) $\Rightarrow \left\{ \begin{array}{l} \text{tmp} = \text{pointer to } \boxed{0} \diagdown \\ \text{tmp}.IntList(1, null); \\ L = \text{tmp}; \end{array} \right.$`

Multiple Constructors and Default Constructors

- All classes have constructors. In the absence of any explicit constructor, get the **default constructor**, as if you had written:

```
public class Foo {  
    public Foo() {}  
}
```

- Multiple **overloaded** constructors are possible, and they can use each other (although the syntax is odd):

```
public class IntList {  
    public IntList(int head, IntList tail) {  
        this.head = head; this.tail = tail;  
    }  
  
    public IntList(int head) {  
        this(head, null); // Calls first constructor.  
    }  
    ...  
}
```

Constructors and Instance Variables

- Instance variables initializations are moved inside constructors that don't start with this(...).

```
class Foo {  
    int x = 5;  
  
    Foo(int y) {  
        DoStuff(y);  
    }  
  
    Foo() {  
        this(42);  
    }  
}  
  
class Foo {  
    int x;  
  
    Foo(int y) {  
        x = 5;  
        DoStuff(y);  
    }  
  
    Foo() {  
        this(42); // Assigns to x  
    }  
}
```

↔

Summary: Java vs. Python

Java	Python
<pre>class Foo { int x = ...; Foo(...) { ... } int f(...) { ... } static int y = 21; static void g(...) { ... } }</pre>	<pre>class Foo: ... x = ... def __init__(self, ...): ... def f(self, ...): ... y = 21 # Referred to as Foo.y @staticmethod def g(...): ...</pre>
<pre>aFoo.f(...) aFoo.x new Foo(...) this</pre>	<pre>aFoo.f(...) aFoo.x Foo(...) self # (typically)</pre>

CS61B Lecture #8: Object-Oriented Mechanisms

Today:

- New in this lecture: the bare mechanics of "object-oriented programming."
- The general topic is: Writing software that operates on many kinds of data.

Overloading

Problem: How to get `System.out.print(x)` to print `x`, regardless of type of `x`?

- In Scheme or Python, one function can take an argument of any type, and then test the type (if needed).
- In Java, methods specify a single type of argument.
- Partial solution: *overloading*—multiple method definitions with the same name and different numbers or types of arguments.
- E.g., `System.out` has type `java.io.PrintStream`, which defines

```
void println() Prints new line.  
void println(String s) Prints S.  
void println(boolean b) Prints "true" or "false"  
void println(char c) Prints single character  
void println(int i) Prints I in decimal  
etc.
```

- Each of these is a different function. Compiler decides which to call on the basis of the arguments' declared types.

Generic Data Structures

Problem: How to get a "list of anything" or "array of anything"?

- Again, no problem in Scheme or Python.
- But in Java, lists (such as `IntList`) and arrays have a single type of element.
- First, the short answer: any *reference* value can be *cast* as (converted to) type `Object` and back, so we can use `Object` as the "generic (reference) type." Such *reference casts* don't change the value of a pointer, but rather tell the compiler how to treat it:

```
Object[] things = new Object[2];
things[0] = new IntList(3, null);
things[1] = "Stuff";
IntList thingsList = (IntList) things[0]; // A cast to IntList
// Both ((IntList) things[0]).head and thingsList.head == 3;
// and ((String) things[1]).startsWith("St") is true, BUT
// things[0].head           Illegal
// things[1].startsWith("St") Illegal
```

And Primitive Values?

- Primitive values (ints, longs, bytes, shorts, floats, doubles, chars, and booleans) are not really convertible to Object.
- Presents a problem for “list of anything.”
- So Java introduced a set of reference types called *wrapper types*, one for each primitive type:

Prim.	Ref.	Prim.	Ref.	Prim.	Ref.
byte	Byte	short	Short	int	Integer
long	Long	char	Character	boolean	Boolean
float	Float	double	Double		

- One can create new wrapper objects for any value (*boxing*):

```
Integer Three = new Integer(3);  
Object ThreeObj = Three;
```

and vice-versa (*unboxing*):

```
int three = Three.intValue();
```

Autoboxing

Boxing and unboxing are automatic (in many cases):

```
Integer Three = 3;  
int three = Three;  
int six = Three + 3;
```

```
Integer[] someInts = { 1, 2, 3 };  
for (int x : someInts) {  
    System.out.println(x);  
}
```

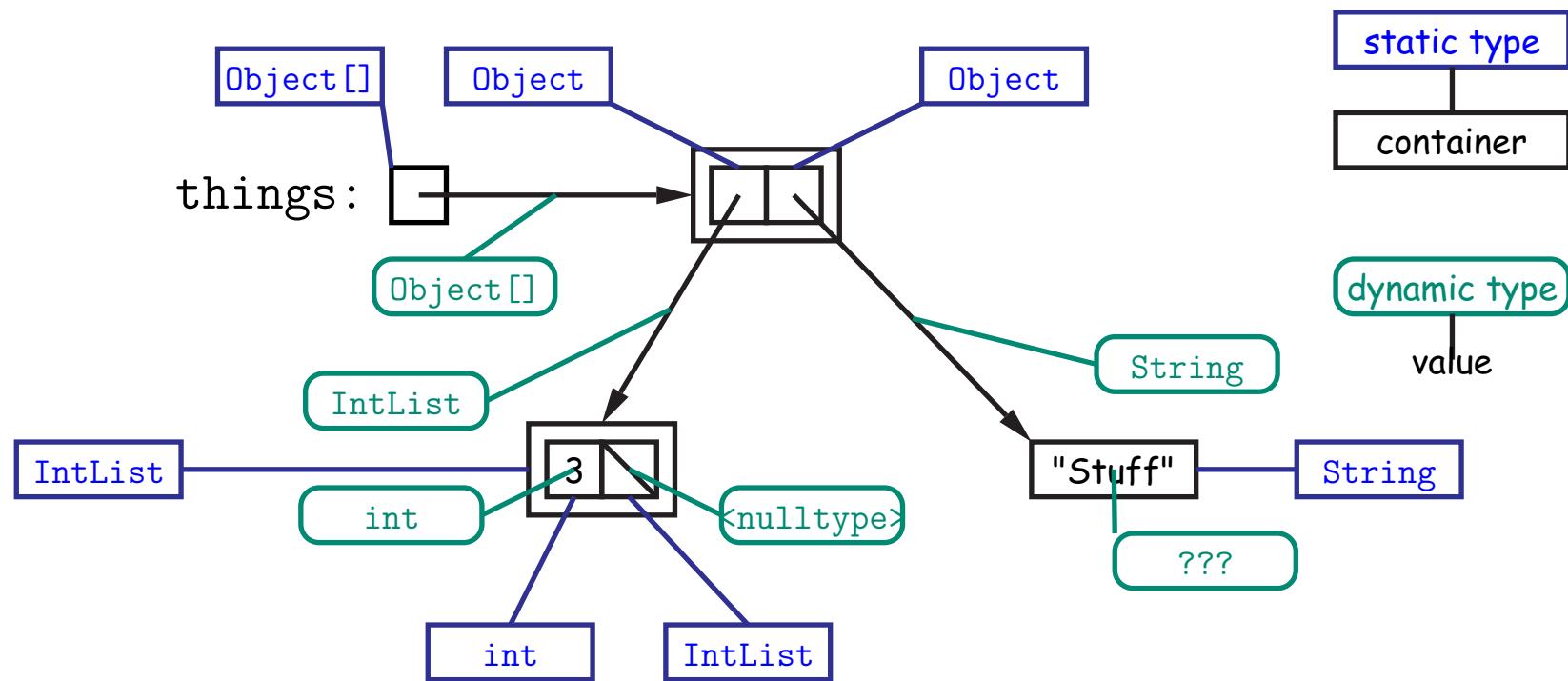
```
System.out.println(someInts[0]);  
// Prints Integer 1, but NOT unboxed.
```

Dynamic vs. Static Types

- Every *value* has a type—its *dynamic type*.
- Every *simple container* (variable, component, parameter), every literal, function call, and operator expression (e.g. $x+y$) has a type—its *static type*.
- Therefore, every *expression* has a static type.
- The static type of an expression is *known to the compiler*.
- The specific dynamic type of the expression—the type of its *value*—is generally unknown to the compiler for reference types.
- But as we shall see, the compiler has some knowledge about the dynamic type of an expression because there is a relationship to the static type of the expression.

Examples of Static and Dynamic Types

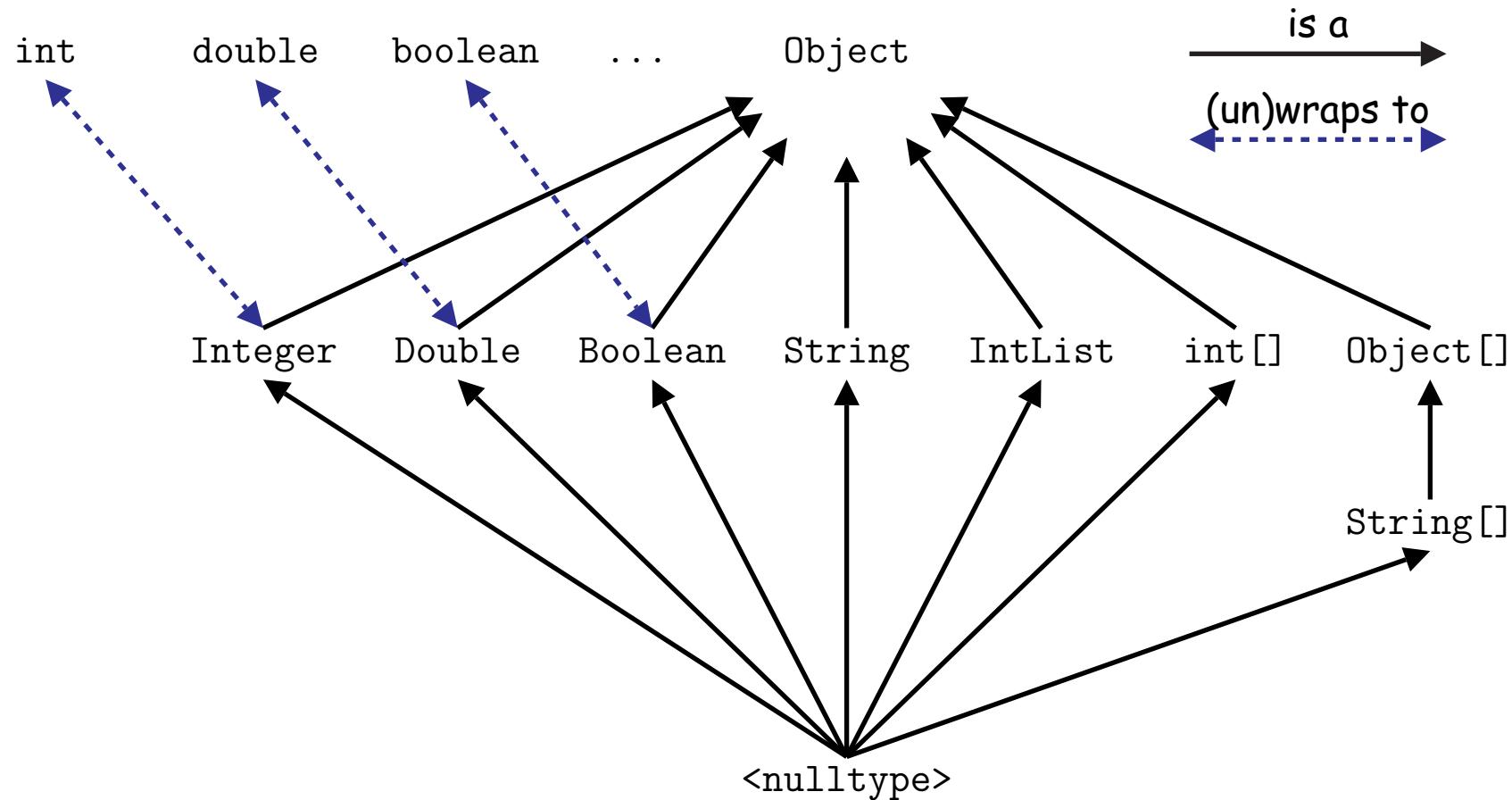
```
Object[] things = new Object[2];
things[0] = new IntList(3, null);
things[1] = "Stuff";
```



Type Hierarchies

- A container with (static) type T may contain a certain value only if that value “is a” T—that is, if the (dynamic) type of the value is a **subtype** of T. Likewise, a function with return type T may return only values that are subtypes of T.
 - All types are subtypes of themselves (& that's all for primitive types)
 - Reference types form a **type hierarchy**; some are subtypes of others.
 - null's type is a subtype of all reference types.
 - All reference types are subtypes of Object.

Java Library Type Hierarchy (Partial)



The Basic Static Type Rule

- Java is designed so that
an expression of (static) type T always yields a value that "is a" T.
- Static types are "known to the compiler," because you declare them, as in

```
String x;           // Static type of field
int f(Object s) { // Static type of call to f, and of parameter
    int y;         // Static type of local variable
```

or they are pre-declared by the language (like 3).

- Compiler insists that in L = E, or in the function call, f(E), where

```
void f(SomeType L) { ... },
```

E's static type must be a subtype of L's static type for reference types.

- Similarly, there are static-type requirements for other operations:
(1) E must have an array type in E[i]; (2) in if (B) ..., B must be boolean; or (3) in E.x, E's static type must define a member named x; etc.

Primitive Types and Coercions

- Primitive types live outside the hierarchy of reference types.
- Although the values of type short, for example, are a subset of those of int, we *don't* say that short is a subtype of int, because they don't quite behave the same.
- However, values of type short can be *coerced* (converted) to a value of type int, using the same cast syntax as for reference types:

```
short x = (short) 3002;
long y = 10000L;
int z = (int) y;
long q = 1000000000000L;
int r = (int) q;
System.out.println(r); // Prints -727379968 (?????)
```

- As the values of r shows, coercions of primitive types, unlike those of reference types, are computations that can change values.

Automatic Coercions, Promotions

- Certain coercions, such converting from short to int, are considered obvious and therefore intrusive.
- So the language silently coerces “smaller” integer types to larger ones, float to double, and integer types to float or double.
- These are called *promotions*.
- Finally, since the compiler can obviously tell what the value of an int literal is, it will convert integer literals to shorter integer types if the values fit:

```
byte x = 127;  
short y = -1024;  
char z = 0x0398; // ⊕
```

Consequences of Compiler's "Sanity Checks"

- These are a *conservative* rules. The last line of the following, which you might think is perfectly sensible, is illegal:

```
int[] A = new int[2];  
Object x = A; // All references are Objects  
A[i] = 0;      // Static type of A is array...  
x[i+1] = 1;    // But not of x: ERROR
```

Compiler figures that not every Object is an array.

- Q: Don't we *know* that x contains array value!?
- A: Yes, but still must tell the compiler, like this:

```
((int[]) x)[i+1] = 1;
```

- Defn: Static type of cast (T) E is T.
- Q: What if x *isn't* an array value, or is null?
- A: For that we have runtime errors—exceptions.

Overriding and Extension

- Notation so far is clumsy.
- Q: If I know Object variable `x` contains a String, why can't I write,
`x.startsWith("this")?`
- A: `startsWith` is only defined on Strings, not on all Objects, so the compiler isn't sure it makes sense, unless you cast.
- But, if an operation *were* defined on all Objects, then you *wouldn't* need clumsy casting.
- Example: `.toString()` is defined on all Objects. You can always say `x.toString()` if `x` has a reference type.
- The default `.toString()` function is not very useful; on an IntList, would produce string like "IntList@2f6684"
- But for any subtype of Object, you may *override* the default definition.

Overriding `toString`

- For example, if `s` is a `String`, `s.toString()` is the identity function (fortunately).
- For any type you define, you may supply your own definition. For example, in `IntList`, could add

```
@Override // Compiler checks that Object really has a toString.  
public String toString() {  
    StringBuffer b = new StringBuffer();  
    b.append("[");  
    for (IntList L = this; L != null; L = L.tail)  
        b.append(" " + L.head);  
    b.append("]");  
    return b.toString();  
}
```

- If `x = new IntList(3, new IntList(4, null))`, then `x.toString()` is "[3 4]".
- Conveniently, various operations requiring `Strings` call `.toString()` for you, so for an `IntList x`, you can write:

```
"Values: " + x, System.out.println(x), or System.out.printf("%s",x);
```

Extending a Class

- To say that class B is a direct subtype of class A (or A is a *direct superclass* of B), write

```
class B extends A { ... }
```

- By default, class ... extends java.lang.Object.
- The subtype *inherits* all fields and methods of its direct superclass (and passes them along to any of its subtypes).
- In class B, you may *override* an instance method (*not* a static method), by providing a new definition with same *signature* (name, return type, argument types).

Rule of Instance Method Calls:

*If f(...) is an instance method, then the call x.f(...) calls whatever overriding of f applies to the *dynamic type* of x, regardless of the static type of x.*

Illustration

```
class Worker {  
    void work() {  
        collectPay();  
    }  
}
```

```
class Prof extends Worker {  
    // Inherits work()  
}  
|
```

```
class TA extends Worker {  
    void work() {  
        while (true) {  
            doLab(); discuss(); officeHour();  
        }  
    }  
}
```

```
Prof paul = new Prof();  
TA daniel = new TA();  
Worker wPaul = paul,  
wDaniel = daniel;  
| paul.work() ==> collectPay();  
| daniel.work() ==> doLab(); discuss(); ...  
| wPaul.work() ==> collectPay();  
| wDaniel.work() ==> doLab(); discuss(); ...  
|
```

Lesson: For instance methods (only), select method based on *dynamic type*. Simple to state, but we'll see it has profound consequences.

What About Fields and Static Methods?

```
class Parent {                                class Child extends Parent {  
    int x = 0;                            String x = "no";  
    static int y = 1;                      static String y = "way";  
    static void f() {                     static void f() {  
        System.out.printf("Ahem!%n");      System.out.printf("I wanna!%n");  
    }                                         }  
    static int f(int x) {                }  
        return x+1;  
    }  
}
```

```
Child tom = new Child(); | tom.x    ==> no      pTom.x    ==> 0  
Parent pTom = tom;     | tom.y    ==> way     pTom.y    ==> 1  
                      | tom.f() ==> I wanna! pTom.f() ==> Ahem!  
                      | tom.f(1) ==> 2       pTom.f(1) ==> 2
```

Lesson: Fields *hide* inherited fields of same name; static methods *hide* methods of the same signature.

Real Lesson: Hiding causes confusion; so understand it, but don't do it!

What's the Point?

- The mechanism described here allows us to define a kind of *generic* method.
- A superclass can define a set of operations (methods) that are common to many different classes.
- Subclasses can then provide different implementations of these common methods, each specialized in some way.
- All subclasses will have at least the methods listed by the superclass.
- So when we write methods that operate on the superclass, they will automatically work for all subclasses with no extra work.

CS61B Lecture #9: Interfaces and Abstract Classes

Recreation

Show that for any polynomial with a leading coefficient of 1 and integral coefficients, all rational roots are integers.

Reminders:

1. The four projects are individual efforts in this class (no partnerships). Feel free to discuss projects or pieces of them before doing the work. But you must complete each project yourself. That is, feel free to discuss projects with each other, but be aware that we expect your work to be substantially different from that of all your classmates (in this or any other semester). You will find a more detailed account of our policy in under the "Course Info" tab on the course website.
2. `git status` is your friend!
3. Please use `git-bug` for submitting problems to us. Piazza and email are really not great for bug reports.

Abstract Methods and Classes

- Instance method can be *abstract*: No body given; must be supplied in subtypes.
- One good use is in specifying a pure interface to a family of types:

```
/** A drawable object. */
public abstract class Drawable {
    // "abstract class" = "can't say new Drawable"
    /** Expand THIS by a factor of XSIZE in the X direction,
     * and YSIZE in the Y direction. */
    public abstract void scale(double xscale, double yscale);

    /** Draw THIS on the standard output. */
    public abstract void draw();
}
```

- Now a `Drawable` is something that has *at least* the operations `scale` and `draw` on it.
- Can't create a `Drawable` because it's abstract.
- In fact, in this case, it wouldn't make any sense to create one, because it has two methods without any implementation.

Methods on Drawables

```
/** A drawable object. */
public abstract class Drawable {
    /** Expand THIS by a factor of XSIZE in the X direction,
     * and YSIZE in the Y direction. */
    public abstract void scale(double xsize, double ysize);
    /** Draw THIS on the standard output. */
    public abstract void draw();
}
```

- Can't write new `Drawable()`, BUT, we can write methods that operate on `Drawables`

```
void drawAll(Drawable[] thingsToDraw) {
    for (Drawable thing : thingsToDraw)
        thing.draw();
}
```

- But `draw` has no implementation! How can this work?

Concrete Subclasses

- Regular classes can extend abstract ones to make them "less abstract" by overriding their abstract methods.
- Can define kinds of `Drawables` that are *concrete* (non-abstract), with all methods having implementations.
- Since all methods are implemented, it makes sense to use `new` on these types—all the method calls make sense.

Concrete Subclass Examples

```
public class Rectangle extends Drawable {  
    public Rectangle(double w, double h) { this.w = w; this.h = h; }  
    public void scale(double xsize, double ysize) {  
        w *= xsize; h *= ysize;  
    }  
    public void draw() { draw a w x h rectangle }  
    private double w,h;  
}
```

Any Oval or Rectangle is a Drawable.

```
public class Oval extends Drawable {  
    public Oval(double xrad, double yrad) {  
        this.xrad = xrad; this.yrad = yrad;  
    }  
    public void scale(double xsize, double ysize) {  
        xrad *= xsize; yrad *= ysize;  
    }  
    public void draw() { draw an oval with axes xrad and yrad }  
    private double xrad, yrad;  
}
```

Using Concrete Classes

- We *can* create new `Rectangles` and `Ovals`.
- Since these classes are subtypes of `Drawable`, we can put them in any container whose static type is `Drawable`, ...
- ... and therefore can pass them to any method that expects `Drawable` parameters:
- Thus, writing

```
Drawable[] things = {  
    new Rectangle(3, 4), new Oval(2, 2)  
};  
drawAll(things);
```

draws a 3×4 rectangle and a circle with radius 2.

Aside: Documentation

- Our style checker would insist on comments for all the methods, constructors, and fields of the concrete subtypes.
- But we already have comments for `draw` and `scale` in the class `Drawable`, and good practice demands that all overridings of these methods should obey at least these comments. Hence, comments are often unhelpful on overriding methods.
- Still, the reader would like to know that a given method `does` override something. Hence, the `@Override` annotation:

```
@Override
public void scale(double xsize, double ysize) {
    xrad *= xsize; yrad *= ysize;
}
@Override
public void draw() { draw a circle with radius rad }
```

- The compiler will check that these method headers are proper overridings of the parent's methods, and our style checker won't complain about the lack of comments.

Interfaces

- In generic English usage, an *interface* is a “point where interaction occurs between two systems, processes, subjects, etc.” (*Concise Oxford Dictionary*).
- In programming, often use the term to mean a *description* of this generic interaction, specifically, a description of the functions or variables by which two things interact.
- Java uses the term to refer to a slight variant of an abstract class that (until Java 1.7) contains only abstract methods (and static constants), like this:

```
public interface Drawable {  
    void scale(double xsize, double ysize); // Automatically public.  
    void draw();  
}
```

- Interfaces are automatically abstract: **can't** say `new Drawable()`; **can** say `new Rectangle(...)`.

Implementing Interfaces

- Idea is to treat Java interfaces as the public **specifications** of data types, and classes as their **implementations**:

```
public class Rectangle implements Drawable { ... }
```

(We **extend** ordinary classes and **implement** interfaces, hence the change in keyword.)

- Can use the interface as for abstract classes:

```
void drawAll(Drawable[] thingsToDo) {
    for (Drawable thing : thingsToDo)
        thing.draw();
}
```

- Again, this works for **Rectangles** and any other implementation of **Drawable**.

Multiple Inheritance

- Can *extend* one class, but *implement* any number of interfaces.
- Contrived Example:

```
interface Readable {  
    Object get();  
}
```

```
interface Writable {  
    void put(Object x);  
}
```

```
class Source implements Readable {  
    public Object get() { ... }  
}
```

```
void copy(Readable r,  
          Writable w) {  
    w.put(r.get());  
}
```

```
class Sink implements Writable {  
    public void put(Object x) { ... }  
}
```

```
class Variable implements Readable, Writable {  
    public Object get() { ... }  
    public void put(Object x) { ... }  
}
```

- The first argument of *copy* can be a *Source* or a *Variable*. The second can be a *Sink* or a *Variable*.

Review: Higher-Order Functions

- In Python, you had *higher-order functions* like this:

```
def map(proc,      items):
    #  function      list
    if items is None:
        return None
    else:
        return IntList(proc(items.head), map(proc, items.tail))
```

and you could write

```
map(abs, makeList(-10, 2, -11, 17))
====> makeList(10, 2, 11, 17)
map(lambda x: x * x, makeList(1, 2, 3, 4))
====> makeList(t(1, 4, 9, 16))
```

- Java does not have these directly, but we can use abstract classes or interfaces and subtyping to get the same effect (with more writing)

Map in Java

```
/** Function with one integer argument */

public interface IntUnaryFunction {
    int apply(int x);
}
```

```
IntList map(IntUnaryFunction proc,
            IntList items) {
    if (items == null)
        return null;
    else return new IntList(
        proc.apply(items.head),
        map(proc, items.tail));
}
```

- It's the use of this function that's clumsy. First, define class for absolute value function; then create an instance:

```
class Abs implements IntUnaryFunction {
    public int apply(int x) { return Math.abs(x); }
}
```

```
R = map(new Abs(), some list);
```

Lambda Expressions

- Since Java 7, one can create classes like Abs on the fly with *anonymous classes*:

```
R = map(new IntUnaryFunction() {
    public int apply(int x) { return Math.abs(x); }
}, some list);
```

- This is sort of like declaring

```
class Anonymous implements IntUnaryFunction {
    public int apply(int x) { return Math.abs(x); }
}
```

and then writing

```
R = map(new Anonymous(), some list);
```

Lambda in Java

- In modern Java, lambda expressions are even more succinct. If a parameter has a type is an interface with a single abstract method (a *functional interface*), Java will figure out the necessary class from just a parameter list and body:

```
R = map((int x) -> Math.abs(x), some list);
```

We don't even need the parameter's type:

```
R = map((x) -> Math.abs(x), some list);
```

and if the body is just a call on a function that already exists:

```
R = map(Math::abs, some list);
```

- These figure out you need an anonymous IntUnaryFunction and create one.
- You can see examples of this sort of thing in `flood.GUI`:

```
addMenuItem("Game->New", this::newGame);
```

Here, the second parameter of `ucb.gui2.TopLevel.addMenuItem` is a *call-back function*.

- It has the Java library type `java.util.function.Consumer`, which has a one-argument method, like `IntUnaryFunction`.

Inheriting Headers vs. Method Bodies

- One can implement multiple interfaces, but extend only one class:
multiple interface inheritance, but single body inheritance.
- This scheme is simple, and pretty easy for language implementors to implement.
- However, there are cases where it would be nice to be able to "mix in" implementations from a number of sources.

Extending Supertypes, Default Implementations

- As indicated above, before Java 8, interfaces contained just static constants and abstract methods.
- Java 8 introduced static methods into interfaces and also *default methods*, which are essentially instance methods and are used whenever a method of a class implementing the interface would otherwise be abstract.
- Suppose I want to add a new one-parameter `scale` method to all concrete subclasses of the interface `Drawable`. Normally, that would involve adding an implementation of that method to all concrete classes.
- We could instead make `Drawable` an abstract class again, but in the general case that can have its own problems.

Default Methods in Interfaces

- So Java 8 introduced default methods:

```
public interface Drawable {  
    void scale(double xsize, double ysize);  
    void draw();  
  
    /** Scale by SIZE in the X and Y dimensions. */  
    default void scale(double size) {  
        scale(size, size);  
    }  
}
```

- Now in any class that implements Drawab;e but does not have a definition of scale with one argument, this method will supply a default implementation for it.
- Useful feature, but, as in other languages with full multiple inheritance (like C++ and Python), it can lead to confusing programs. I suggest you use them sparingly.

Fancy Example From the Java Library

```
package java.util.function;

public interface Consumer<T> {
    void accept(T t);

    default Consumer<T> andThen(Consumer<? super T> after) {
        return (x) -> { accept(x); after.accept(x); }
    }
}
```

- Don't worry about the weird stuff in `< >` for now. We'll get to that when we look at *generic definitions*.
- The `andThen` method is another example of Java's version of higher-order function:

```
Consumer<String> print1 = (x) -> System.out.print(x);
print1.accept("Hello");      // Prints "Hello"
Consumer<String> print2 = print1.andThen((y) -> System.out.print(y));
print2.accept("Hello");      // Prints "HelloHello"
```

CS61B Lecture #10: OOP Mechanisms and Class Design

Review: A Puzzle

```
class A {  
    void f() {  
        System.out.println("A.f");  
    }  
    void g() { f(); /* or this.f() */ }  
}
```

```
class B extends A {  
    void f() {  
        System.out.println("B.f");  
    }  
}
```

```
class C {  
    static void main(String[] args) {  
        B aB = new B();  
        h(aB);  
    }  
  
    static void h(A x) { x.g(); }  
}
```

1. What is printed?
2. If we made g static?
3. If we made f static?
4. If we overrode g in B?
5. If f not defined in A?

Choices

- a. A.f
- b. B.f
- c. Some kind of error

Review: A Puzzle

```
class A {  
    void f() {  
        System.out.println("A.f");  
    }  
    void g() { f(); /* or this.f() */ }  
}
```

```
class B extends A {  
    void f() {  
        System.out.println("B.f");  
    }  
}
```

```
class C {  
    static void main(String[] args) {  
        B aB = new B();  
        h(aB);  
    }  
  
    static void h(A x) { x.g(); }  
}
```

1. What is printed?
2. If we made g static?
3. If we made f static?
4. If we overrode g in B?
5. If f not defined in A?

Choices

- a. A.f
- b. B.f
- c. Some kind of error

Review: A Puzzle

```
class A {  
    void f() {  
        System.out.println("A.f");  
    }  
    static void g(A y) { y.f(); }  
}
```

```
class B extends A {  
    void f() {  
        System.out.println("B.f");  
    }  
}
```

```
class C {  
    static void main(String[] args) {  
        B aB = new B();  
        h(aB);  
    }  
}
```

```
static void h(A x) { A.g(x); } // x.g(x) also
```

legal here

```
}
```

1. What is printed?
2. If we made g static?
3. If we made f static?
4. If we overrode g in B?
5. If f not defined in A?

Choices

- a. A.f
- b. B.f
- c. Some kind of error

Review: A Puzzle

```
class A {  
    void f() {  
        System.out.println("A.f");  
    }  
    static void g(A y) { y.f(); }  
}
```

```
class B extends A {  
    void f() {  
        System.out.println("B.f");  
    }  
}
```

```
class C {  
    static void main(String[] args) {  
        B aB = new B();  
        h(aB);  
    }  
}
```

```
static void h(A x) { A.g(x); } // x.g(x) also
```

legal here

```
}
```

1. What is printed?
2. If we made g static?
3. If we made f static?
4. If we overrode g in B?
5. If f not defined in A?

Choices

- a. A.f
- b. B.f
- c. Some kind of error

Review: A Puzzle

```
class A {  
    static void f() {  
        System.out.println("A.f");  
    }  
    void g() { f(); /* or this.f() */ }  
}
```

```
class B extends A {  
    static void f() {  
        System.out.println("B.f");  
    }  
}
```

```
class C {  
    static void main(String[] args) {  
        B aB = new B();  
        h(aB);  
    }  
  
    static void h(A x) { x.g(); }  
}
```

1. What is printed?
2. If we made g static?
3. If we made f static?
4. If we overrode g in B?
5. If f not defined in A?

Choices

- a. A.f
- b. B.f
- c. Some kind of error

Review: A Puzzle

```
class A {  
    static void f() {  
        System.out.println("A.f");  
    }  
    void g() { f(); /* or this.f() */ }  
}
```

```
class B extends A {  
    static void f() {  
        System.out.println("B.f");  
    }  
}
```

```
class C {  
    static void main(String[] args) {  
        B aB = new B();  
        h(aB);  
    }  
  
    static void h(A x) { x.g(); }  
}
```

1. What is printed?
2. If we made g static?
3. If we made f static?
4. If we overrode g in B?
5. If f not defined in A?

Choices

- a. A.f
- b. B.f
- c. Some kind of error

Review: A Puzzle

```
class A {  
    void f() {  
        System.out.println("A.f");  
    }  
    void g() { f(); /* or this.f() */ }  
}
```

```
class B extends A {  
    void f() {  
        System.out.println("B.f");  
    }  
    void g() { f(); }  
}
```

```
class C {  
    static void main(String[] args) {  
        B aB = new B();  
        h(aB);  
    }  
  
    static void h(A x) { x.g(); }  
}
```

1. What is printed?
2. If we made g static?
3. If we made f static?
4. If we overrode g in B?
5. If f not defined in A?

Choices

- a. A.f
- b. B.f
- c. Some kind of error

Review: A Puzzle

```
class A {  
    void f() {  
        System.out.println("A.f");  
    }  
    void g() { f(); /* or this.f() */ }  
}
```

```
class B extends A {  
    void f() {  
        System.out.println("B.f");  
    }  
    void g() { f(); }  
}
```

```
class C {  
    static void main(String[] args) {  
        B aB = new B();  
        h(aB);  
    }  
  
    static void h(A x) { x.g(); }  
}
```

1. What is printed?
2. If we made g static?
3. If we made f static?
4. If we overrode g in B?
5. If f not defined in A?

Choices

- a. A.f
- b. B.f
- c. Some kind of error

Review: A Puzzle

```
class A {  
    void g() { f(); /* or this.f() */ }  
}
```

```
class B extends A {  
    void f() {  
        System.out.println("B.f");  
    }  
}
```

```
class C {  
    static void main(String[] args) {  
        B aB = new B();  
        h(aB);  
    }  
  
    static void h(A x) { x.g(); }  
}
```

1. What is printed?
2. If we made g static?
3. If we made f static?
4. If we overrode g in B?
5. If f not defined in A?

Choices

- a. A.f
- b. B.f
- c. Some kind of error

Review: A Puzzle

```
class A {  
    void g() { f(); /* or this.f() */ }  
}
```

```
class B extends A {  
    void f() {  
        System.out.println("B.f");  
    }  
}
```

```
class C {  
    static void main(String[] args) {  
        B aB = new B();  
        h(aB);  
    }  
  
    static void h(A x) { x.g(); }  
}
```

1. What is printed?
2. If we made g static?
3. If we made f static?
4. If we overrode g in B?
5. If f not defined in A?

Choices

- a. A.f
- b. B.f
- c. Some kind of error

Answer to Puzzle

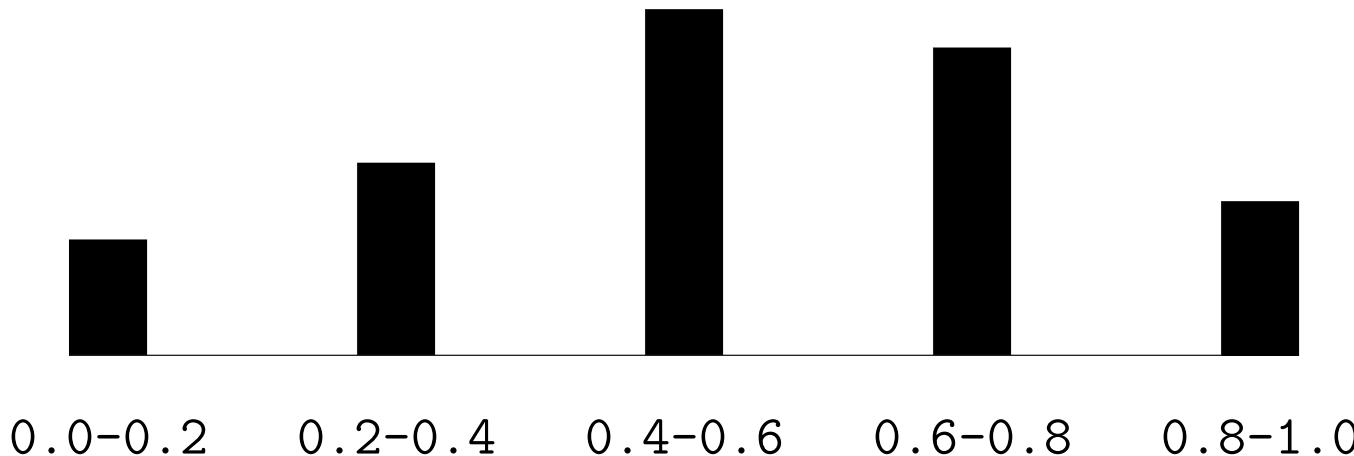
1. Executing `java C` prints _____, because
 - A. `C.main` calls `h` and passes it `aB`, whose dynamic type is `B`.
 - B. `h` calls `x.g()`. Since `g` is inherited by `B`, we execute the code for `g` in class `A`.
 - C. `g` calls `this.f()`. Now `this` contains the value of `h`'s argument, whose dynamic type is `B`. Therefore, we execute the definition of `f` that is in `B`.
 - D. In calls to `f`, in other words, static type is ignored in figuring out what method to call.
2. If `g` were static, we see _____; selection of `f` still depends on dynamic type of `this`. Same for overriding `g` in `B`.
3. If `f` were static, would print _____ because then selection of `f` would depend on static type of `this`, which is `A`.
4. If `f` were not defined in `A`, we'd see _____

Answer to Puzzle

1. Executing `java C` prints B.f, because
 - A. C.main calls h and passes it aB, whose dynamic type is B.
 - B. h calls x.g(). Since g is inherited by B, we execute the code for g in class A.
 - C. g calls this.f(). Now this contains the value of h's argument, whose dynamic type is B. Therefore, we execute the definition of f that is in B.
 - D. In calls to f, in other words, static type is ignored in figuring out what method to call.
2. If g were static, we see B.f; selection of f still depends on dynamic type of this. Same for overriding g in B.
3. If f were static, would print A.f because then selection of f would depend on static type of this, which is A.
4. If f were not defined in A, we'd see a compile-time error

Example: Designing a Class

Problem: Want a class that represents histograms, like this one:



Analysis: What do we need from it? At least:

- Specify buckets and limits.
- Accumulate counts of values.
- Retrieve counts of values.
- Retrieve numbers of buckets and other initial parameters.

Specification Seen by Clients

- The *clients* of a module (class, program, etc.) are the programs or methods that *use* that module's exported definitions.
- In Java, intention is that exported definitions are designated **public**.
- Clients are intended to rely on *specifications*, (aka APIs) not code.
- *Syntactic specification*: method and constructor headers—syntax needed to use.
- *Semantic specification*: what they do. No formal notation, so use comments.
 - Semantic specification is a *contract*.
 - Conditions client must satisfy (*preconditions*, marked "Pre:" in examples below).
 - Promised results (*postconditions*).
 - Design these to be *all the client needs!*
 - Exceptions communicate errors, specifically failure to meet pre-conditions.

Histogram Specification and Use

```
/** A histogram of floating-point values */
public interface Histogram {
    /** The number of buckets in THIS. */
    int size();

    /** Lower bound of bucket #K. Pre: 0<=K<size(). */
    double low(int k);

    /** # of values in bucket #K. Pre: 0<=K<size(). */
    int count(int k);

    /** Add VAL to the histogram. */
    void add(double val);
}
```

Sample output:

```
>= 0.00 | 10
>= 10.25 | 80
>= 20.50 | 120
>= 30.75 | 50
```

```
void fillHistogram(Histogram H,
                   Scanner in) {
    while (in.hasNextDouble())
        H.add(in.nextDouble());
}

void printHistogram(Histogram H) {
    for (int i = 0; i < H.size(); i += 1)
        System.out.printf
            (">=%5.2f | %4d%n",
             H.low(i), H.count(i));
}
```

An Implementation

```
public class FixedHistogram implements Histogram {  
    private double _low, _high; /* From constructor*/  
    private int[] _count; /* Value counts */  
  
    /** A new histogram with SIZE buckets of values >= LOW and < HIGH. */  
    public FixedHistogram(int size, double low, double high)  
    {  
        if (low >= high || size <= 0) throw new IllegalArgumentException();  
        _low = low; _high = high;  
        _count = new int[size];  
    }  
  
    public int size() { return _count.length; }  
    public double low(int k) { return _low + k*(_high-_low)/_count.length; }  
  
    public int count(int k) { return _count[k]; }  
  
    public void add(double val) {  
        if (val >= _low && val < _high)  
            _count[(int) ((val-_low)/(_high-_low) * _count.length)] += 1;  
    }  
}
```

Let's Make a Tiny Change

Don't require *a priori* bounds:

```
class FlexHistogram implements Histogram {  
    /** A new histogram with SIZE buckets. */  
    public FlexHistogram(int size) {  
        ?  
    }  
    // What needs to change?  
}
```

- How would you do this? Profoundly changes implementation.
- But *clients* (like `printHistogram` and `fillHistogram`) still work with no changes.
- Illustrates the power of *separation of concerns*.

Implementing the Tiny Change

- Pointless to pre-allocate the `_count` array.
- Don't know bounds, so must save arguments to `add`.
- Then recompute `_count` array "lazily" when `_count(...)` called.
- Invalidate `_count` array whenever histogram changes.

```
class FlexHistogram implements Histogram {  
    private ArrayList<Double> _values = new ArrayList<>();  
    int _size;  
    private int[] _count;  
  
    public FlexHistogram(int size) { _size = size; _count = null; }  
  
    public int size() { return _size; }  
  
    public void add(double x) { _count = null; _values.add(x); }  
  
    public int count(int k) {  
        if (_count == null) { compute _count from _values here. }  
        return _count[k];  
    }  
}
```

Advantages of Procedural Interface over Visible Fields

By using public method for `count` instead of making the array `_count` visible, the “tiny change” is transparent to clients:

- If client had to write `myHist._count[k]`, it would mean
“The number of items currently in the k^{th} bucket of histogram `myHist` (which, by the way, is stored in an array called `_count` in `myHist` that always holds the up-to-date count).”
- Parenthetical comment *worse than useless* to the client.
- If `_count` array had been visible, after “tiny change,” *every use* of `count` in client program would have to change.
- So using a method for the public `count` method decreases what client *has to* know, and (therefore) has to change.

Recreation

An integer is divided by 9 when a certain one of its digits is deleted, and the resulting number is again divisible by 9.

- a. Prove that actually dividing the resulting number by 9 results in deleting another digit.
- b. Find all integers satisfying the conditions of this problem.

CS61B Lecture #11: Examples: Comparable & Reader + Some Features Supporting Abstraction

Comparable

- Java library provides an interface to describe Objects that have a *natural order* on them, such as `String`, `Integer`, `BigInteger` and `BigDecimal`:

```
public interface Comparable { // For now, the Java 1.4 version
    /** Returns value <0, == 0, or > 0 depending on whether THIS is
     * <, ==, or > OBJ. Exception if OBJ not of compatible type. */
    int compareTo(Object obj);
}
```

- Might use in a general-purpose max function:

```
/** The largest value in array A, or null if A empty. */
public static Comparable max(Comparable[] A) {
    if (A.length == 0) return null;
    Comparable result; result = A[0];
    for (int i = 1; i < A.length; i += 1)
        if (result.compareTo(A[i]) < 0) result = A[i];
    return result;
}
```

- Now `max(S)` will return maximum value in S if S is an array of `Strings`, or any other kind of Object that implements `Comparable`.

Examples: Implementing Comparable

```
/** A class representing a sequence of ints. */
class IntSequence implements Comparable {
    private int[] _myValues;
    private int _myCount;
    ...
    public int get(int k) { return _myValues[k]; }

    @Override
    public int compareTo(Object obj) {
        IntSequence x = (IntSequence) obj; // Blows up if obj not an IntSequence
        for (int i = 0; i < _myCount && i < x._myCount; i += 1) {
            if (_myValues[i] < x._myValues[i]) {
                return -1;
            } else if (_myValues[i] > x._myValues[i]) {
                return 1;
            }
        }
        return _myCount - x._myCount; // <0 iff _myCount < x._myCount
    }
}
```

Implementing Comparable II

- Also possible to add an interface retroactively.
- If `IntSequence` did not implement `Comparable`, but did implement `compareTo` (without `@Override`), we could write

```
class ComparableIntSequence extends IntSequence implements Comparable {  
}
```

- Java would then "match up" the `compareTo` in `IntSequence` with that in `Comparable`.

Java Generics (I)

- We've shown you the old Java 1.4 Comparable. The current version uses a newer feature: Java generic types:

```
public interface Comparable<T> {
    int compareTo(T x);
}
```

- Here, T is like a formal parameter in a method, except that its "value" is a type.
- Revised IntSequence (no casting needed):

```
class IntSequence implements Comparable<IntSequence> {
    ...
    @Override
    public int compareTo(IntSequence x) {
        for (int i = 0; i < _myCount && i < x._myCount; i += 1) {
            if (_myValues[i] < x._myValues[i]) ...

            return _myCount - x._myCount;
        }
    }
}
```

Example: Readers

- Java class `java.io.Reader` abstracts *sources of characters*.
- Here, we present a revisionist version (not the real thing):

```
public interface Reader { // Real java.io.Reader is abstract class
    /** Release this stream: further reads are illegal */
    void close();

    /** Read as many characters as possible, up to LEN,
     *  into BUF[OFF], BUF[OFF+1],..., and return the
     *  number read, or -1 if at end-of-stream. */
    int read(char[] buf, int off, int len);

    /** Short for read(BUF, 0, BUF.length). */
    int read(char[] buf);

    /** Read and return single character, or -1 at end-of-stream. */
    int read();
}
```

- Can't write new `Reader()`; it's abstract. So what good is it?

Generic Partial Implementation

- According to their specifications, some of Reader's methods are related.
- Can express this with a *partial implementation*, which leaves key methods unimplemented and provides default bodies for others.
- Result still abstract: can't use new on it.

```
/** A partial implementation of Reader. Concrete
 * implementations MUST override close and read(,,).
 * They MAY override the other read methods for speed. */
public abstract class AbstractReader implements Reader {
    // Next two lines are redundant.
    public abstract void close();
    public abstract int read(char[] buf, int off, int len);

    public int read(char[] buf) { return read(buf,0,buf.length); }

    public int read() { return (read(_buf1) == -1) ? -1 : _buf1[0]; }

    private char[] _buf1 = new char[1];
}
```

Implementation of Reader: `StringReader`

The class `StringReader` reads characters from a `String`:

```
public class StringReader extends AbstractReader {  
    private String _str;  
    private int _k;  
    /** A Reader that delivers the characters in S. */  
    public StringReader(String s) {  
        _str = s; _k = 0;  
    }  
  
    public void close() { _str = null; }  
  
    public int read(char[] buf, int off, int len) {  
        if (_k == _str.length())  
            return -1;  
        len = Math.min(len, _str.length() - _k);  
        _str.getChars(_k, _k + len, buf, off);  
        _k += len;  
        return len;  
    }  
}
```

Using Reader

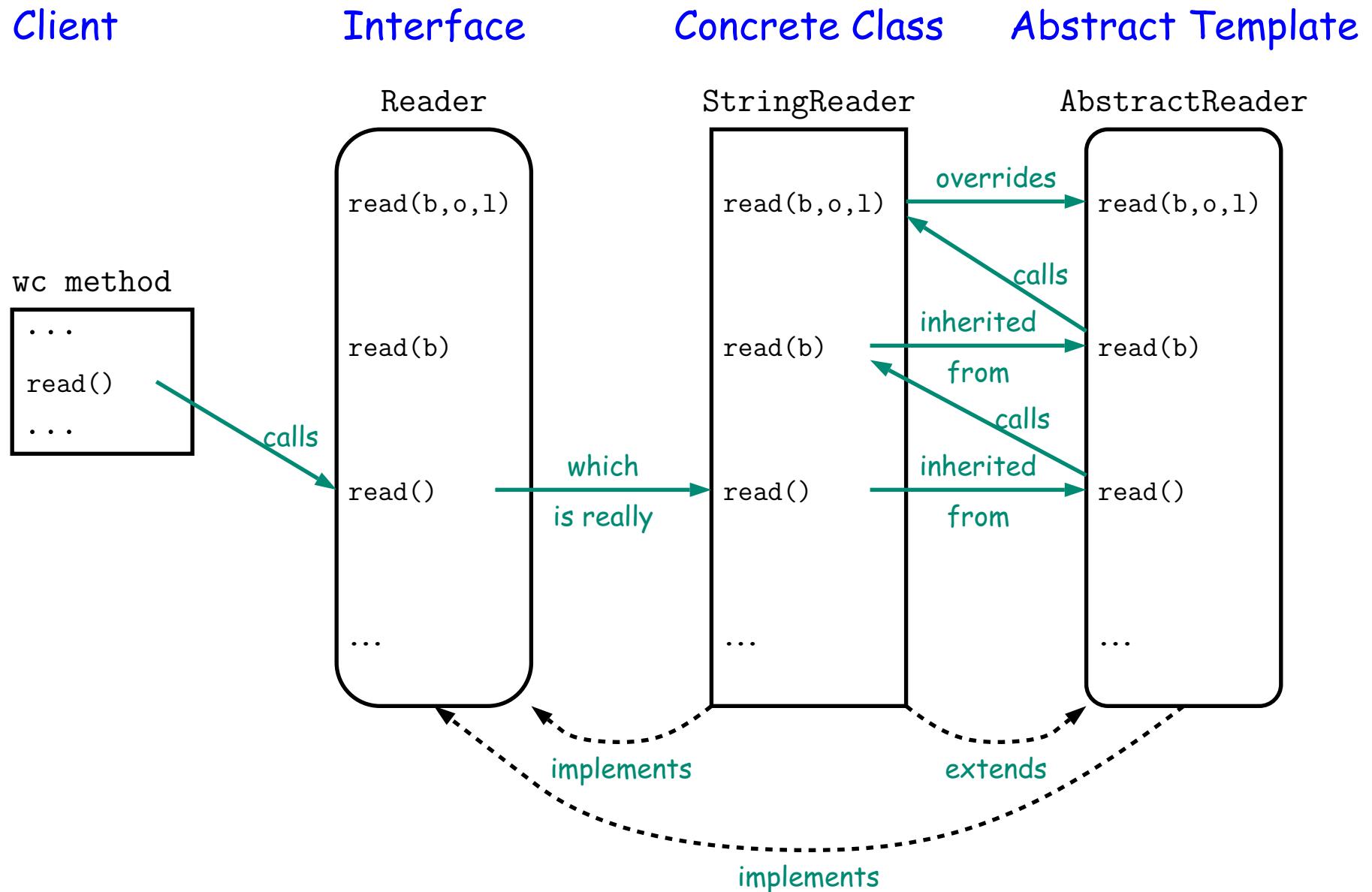
Consider this method, which counts words:

```
/** The total number of words in R, where a "word" is
 * a maximal sequence of non-whitespace characters. */
int wc(Reader r) {
    int c0, count;
    c0 = ' '; count = 0;
    while (true) {
        int c = r.read();
        if (c == -1) return count;
        if (Character.isWhitespace((char) c0)
            && !Character.isWhitespace((char) c))
            count += 1;
        c0 = c;
    }
}
```

This method works for *any* Reader:

```
wc(new StringReader(someText))          // # words in someText
wc(new InputStreamReader(System.in))    // # words in standard input
wc(new FileReader("foo.txt"))           // # words in file foo.txt.
```

How It Fits Together



Lessons

- The `Reader` interface class served as a *specification* for a whole set of readers.
- Ideally, most client methods that deal with `Readers`, like `wc`, will specify type `Reader` for the formal parameters, not a specific kind of `Reader`, thus assuming as little as possible.
- And only when a client creates a new `Reader` will it get specific about what subtype of `Reader` it needs.
- That way, client's methods are as *widely applicable* as possible.
- Finally, `AbstractReader` is a tool for implementors of non-abstract `Reader` classes, and not used by clients.
- Alas, Java library is not pure. E.g., `AbstractReader` is really just called `Reader` and there is no interface. In this example, we saw what they *should* have done!
- The `Comparable` interface allows definition of functions that depend only on a limited subset of the properties (methods) of their arguments (such as "must have a `compareTo` method").

More OOP Features Supporting Abstraction

Parent Constructors

- In lecture notes #7, talked about how Java allows implementer of a class to control all manipulation of objects of that class.
- In particular, this means that Java gives the constructor of a class the first shot at each new object.
- When one class extends another, there are two constructors—one for the parent type and one for the new (child) type.
- In this case, Java guarantees that *one of the parent's constructors is called first*. In effect, there is a call to a parent constructor at the beginning of every one of the child's constructors.
- You can call the parent's constructor yourself explicitly.

```
class Figure {  
    public Figure(int sides) {  
        ...  
    }...  
}
```

```
class Rectangle extends Figure {  
    public Rectangle() {  
        super(4);  
    }...  
}
```

Default Constructors

- By default, Java calls the “default” (parameterless) constructor if there is no explicit constructor called.

```
/* This... */          /* Is equivalent to... */  
class Thingy extends Rectangle {    class Thingy extends Rectangle {  
    public Thingy() {        public Thingy() {  
        setThingsUp();        super();  
    }                    setThingsUp();  
}                      }  
}                      }
```

- And it creates a default constructor for a class if no other constructor is defined for the class.

```
/* This... */      /* Is equivalent to... */      /* And thus to... */  
class Crate {      class Crate {                class Crate {  
}                      public Crate() {        public Crate() {  
}                      }                    super();  
}                      }  
}                      }
```

What Happens Here?

```
class Figure {  
    public Figure(int sides) {  
        ...  
    }  
}  
  
class Rectangle extends Figure {  
}
```

What Happens Here?

```
class Figure {  
    public Figure(int sides) {  
        ...  
    }  
}  
  
class Rectangle extends Figure {  
}
```

Answer: Compiler error. Rectangle has an implicit constructor that tries to call the default constructor in Figure, but there isn't one.

Using an Overridden Method

- Suppose that you wish to *add* to the action defined by a superclass's method, rather than to completely override it.
- The overriding method can refer to overridden methods by using the special prefix `super`.
- For example, you have a class with expensive functions, and you'd like a memoizing version of the class.

```
class ComputeHard {  
    int cogitate(String x, int y) { ... }  
}  
  
class ComputeLazily extends ComputeHard {  
    int cogitate(String x, int y) {  
        if (don't already have answer for this x and y) {  
            int result = super.cogitate(x, y); // <<< Calls overridden function  
            memoize (save) result;  
            return result;  
        }  
        return memoized result;  
    }  
}
```

CS61B Lecture #12: Delegation, Exceptions, Assorted Features

- Delegation
- Exceptions
- Importing
- Nested classes.
- Type testing.

Trick: Delegation and Wrappers

- Not always appropriate to use inheritance to extend something.
- Homework gives example of a TrReader, which *contains* another Reader, to which it *delegates* the task of actually going out and reading characters.
- Another example: a class that instruments objects:

```
interface Storage {  
    void put(Object x);  
    Object get();  
}
```

}

```
// ORIGINAL  
Storage S = something;  
f(S);
```

```
class Monitor implements Storage {  
    int gets, puts;  
    private Storage store;  
    Monitor(Storage x) { store = x; gets = puts = 0; }  
    public void put(Object x) { puts += 1; store.put(x); }  
    public Object get() { gets += 1; return store.get(); }  
}
```

```
// INSTRUMENTED  
Monitor S = new Monitor(something);  
f(S);  
System.out.println(S.gets + " gets");
```

Monitor is called a *wrapper class*.

What to do About Errors?

- Large amount of any production program devoted to detecting and responding to errors.
- Some errors are external (bad input, network failures); others are internal errors in programs.
- When method has stated precondition, it's the client's job to comply.
- Still, it's nice to detect and report client's errors.
- In Java, we *throw exception objects*, typically:

```
throw new SomeException(optional description);
```

- Exceptions are objects. By convention, they are given two constructors: one with no arguments, and one with a descriptive string argument (which the exception stores).
- Java system throws some exceptions implicitly, as when you dereference a null pointer, or exceed an array bound.

Catching Exceptions

- A **throw** causes each active method call to *terminate abruptly*, until (and unless) we come to a **try** block.
- Catch exceptions and do something corrective with **try**:

```
try {  
    Stuff that might throw exception;  
} catch (SomeException e) {  
    Do something reasonable;  
} catch (SomeOtherException e) {  
    Do something else reasonable;  
}  
Go on with life;
```

- When *SomeException* exception occurs during "Stuff..." and is not handled there, we immediately "do something reasonable" and then "go on with life."
- Descriptive string (if any) available as `e.getMessage()` for error messages and the like.

Catching Exceptions, II

- Using a supertype as the parameter type in a **catch** clause will catch any subtype of that exception as well:

```
try {
    Code that might throw a FileNotFoundException or a
    MalformedURLException ;
} catch (IOException ex) {
    Handle any kind of IOException;
}
```

- Since `FileNotFoundException` and `MalformedURLException` both inherit from `IOException`, the **catch** handles both cases.
- Subtyping means that multiple **catch** clauses can apply; Java takes the first.
- Stylistically, it's nice to be more specific (concrete) about exception types where possible.
- In particular, our style checker will therefore balk at the use of `Exception`, `RuntimeException`, `Error`, and `Throwable` as exception supertypes.

Catching Exceptions, III

- There's a relatively new shorthand for handling multiple exceptions the same way:

```
try {  
    Code that might throw IllegalArgumentException  
    or IllegalStateException;  
} catch (IllegalArgumentException|IllegalStateException ex) {  
    Handle exception;  
}
```

Exceptions: Checked vs. Unchecked

- The object thrown by `throw` command must be a subtype of `Throwable` (in `java.lang`).
- Java pre-declares several such subtypes, among them
 - `Error`, used for serious, unrecoverable errors;
 - `Exception`, intended for all other exceptions;
 - `RuntimeException`, a subtype of `Exception` intended mostly for programming errors too common to be worth declaring.
- Pre-declared exceptions are all subtypes of one of these.
- Any subtype of `Error` or `RuntimeException` is said to be **unchecked**.
- All other exception types are **checked**.

Unchecked Exceptions

- Intended for
 - Programmer errors: many library functions throw `IllegalArgumentException` when one fails to meet a precondition.
 - Errors detected by the basic Java system: e.g.,
 - * Executing `x.y` when `x` is null,
 - * Executing `A[i]` when `i` is out of bounds,
 - * Executing `(String) x` when `x` turns out not to point to a String.
 - Certain catastrophic failures, such as running out of memory.
- May be thrown anywhere at any time with no special preparation.

Checked Exceptions

- Intended to indicate exceptional circumstances that are *expected* to happen from time to time. Examples:
 - Attempting to open a file that does not exist.
 - Input or output errors on a file.
 - Receiving an interrupt.
- Every checked exception that can occur inside a method must either be handled by a try statement, or reported in the method's declaration.
- For example,

```
void myRead() throws IOException, InterruptedException { ... }
```

means that myRead (or something it calls) *might* throw IOException or InterruptedException.

A Language Design Issue

Java makes the following illegal for checked exceptions like IOException.
Why?

```
class Parent {           class Child extends Parent {  
    void f() { ... }     void f() throws IOException { ... }  
}
```

A Language Design Issue

Java makes the following illegal for checked exceptions like IOException.
Why?

```
class Parent {           class Child extends Parent {  
    void f() { ... }      void f() throws IOException { ... }  
}
```

Consider, for example,

```
static void process(Parent p) {  
    p.f();  
}
```

According to the specification for class Parent, this is supposed to be OK, but the call p.f() actually calls Child.f, which might throw IOException. So contrary to the intent of checked exceptions, the process method might throw a checked exception that it does not list.

Good Practice

- Throw exceptions rather than using print statements and System.exit everywhere,
- ... because response to a problem may depend on the *caller*, not just method where problem arises.
- Nice to throw an exception when programmer violates preconditions.
- Particularly good idea to throw an exception rather than let bad input corrupt a data structure.
- Good idea to document when methods throw exceptions.
- To convey information about the cause of exceptional condition, put it into the exception rather than into some global variable:

```
class MyBad extends Exception {  
    public IntList errs;  
    MyBad(IntList nums) { errs=nums; }  
}  
  
try {...  
} catch (MyBad e) {  
    ... e.errs ...  
}
```

Terminology

- Many students speak of "throwing" an error when they mean "throwing an exception in order to *report* an error."
- This is a confusion of implementation (the code a program uses to *internally* signal an exceptional condition) with visible behavior (printing an error message).
- Users are not supposed to see them, but rather the messages that interpret these exceptions and report the exceptional conditions:

Good: Program prints

File myData.txt not found.

Bad: Program prints

Exception in thread "main" java.io.IOException: File not found
at foo.main(foo.java:4)

Importing

- Writing `java.util.List` every time you mean `List` or `java.lang.regex.Pattern` every time you mean `Pattern` is annoying.
- The purpose of the `import` clause at the beginning of a source file is to define abbreviations:
 - `import java.util.List;` means “within this file, you can use `List` as an abbreviation for `java.util.List`.
 - `import java.util.*;` means “within this file, you can use *any* class name in the package `java.util` without mentioning the package.”
- Importing does *not* grant any special access; it *only* allows abbreviation.
- In effect, your program always contains `import java.lang.*;`

Static importing

- One can easily get tired of writing `System.out` and `Math.sqrt`. Do you really need to be reminded with each use that `out` is in the `java.lang.System` package and that `sqrt` is in the `Math` package (duh)?
- Both examples are of *static* members. A feature of Java allows you to abbreviate such references:
 - `import static java.lang.System.out;` means “within this file, you can use `out` as an abbreviation for `System.out`.
 - `import static java.lang.System.*;` means “within this file, you can use *any* static member name in `System` without mentioning the package.
- Again, this is *only* an abbreviation. No special access.
- Alas, you can't do this for classes in the anonymous package.

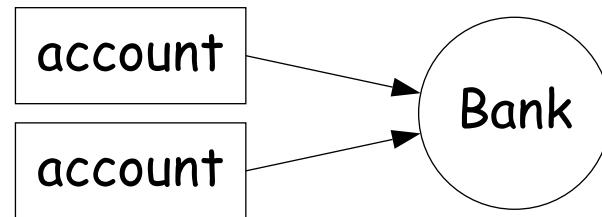
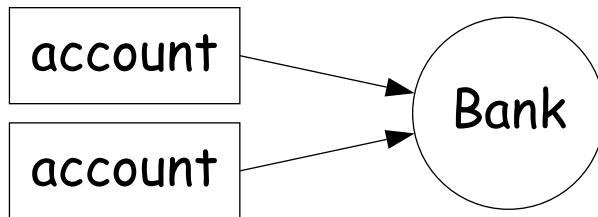
Nesting Classes

- Sometimes, it makes sense to *nest* one class in another. The nested class might
 - be used only in the implementation of the other, or
 - be conceptually "subservient" to the other
- Nesting such classes can help avoid name clashes or "pollution of the name space" with names that will never be used anywhere else.
- Example: Polynomials can be thought of as sequences of terms. Assuming that terms aren't used outside of Polynomials, you might define a class to represent a term *inside* the Polynomial class:

```
class Polynomial {  
  
    methods on polynomials  
  
    private Term[] terms;  
    private static class Term {  
        ...  
    }  
}
```

Inner Classes

- Last slide showed a *static* nested class. Static nested classes are just like any other, except that they can be private or protected, and they can see private variables of the enclosing class.
- Non-static nested classes are called *inner classes*.
- Used when each instance of the nested class is created by and naturally associated with an instance of the containing class, like Banks and Accounts:



Example: Banks and Accounts

```
class Bank {  
    private void addFunds(int amount) { totalFunds += amount; }  
  
    public class Account {  
        int _balance;  
        public void deposit(int amount) {  
            _balance += amount;  
            Bank.this.addFunds(_balance);  
        } // Bank.this means "the bank that created me"  
    }  
}
```

```
Bank bank = new Bank(...);  
Bank.Account a0 = bank.new Account(...);  
Bank.Account a1 = bank.new Account(...);
```

Example: Iterators

```
public class ArrayList<T> {  
    ...  
    public T get(int k) { ... }  
    public int size() { ... }  
  
    public Iterator<T> iterator() {  
        return new ArrayIterator<T>();  
        // or return this.new ArrayIterator<T>();  
    }  
  
    private class ArrayIterator implements Iterator<T> {  
        int _k;  
        ArrayIterator() { _k = 0; }  
        public boolean hasNext() { return _k < size(); }  
        public T next() { _k += 1; return get(_k - 1); }  
    }  
}
```

Type testing: instanceof

- It is possible to ask about the dynamic type of something:

```
void typeChecker(Reader r) {  
    if (r instanceof TrReader)  
        System.out.print("Translated characters: ");  
    else  
        System.out.print("Characters: ");  
    ...  
}
```

- However, this is seldom what you want to do. Why do this:

```
if (x instanceof StringReader)  
    read from (StringReader) x;  
else if (x instanceof FileReader)  
    read from (FileReader) x;  
...
```

when you can just call x.read()?!

- In general, use instance methods rather than instanceof.

CS61B Lecture #14: Packages, Access

Package Mechanics

- Classes correspond to things being modeled (represented) in one's program.
- Packages are collections of "related" classes and other packages.
- Java puts standard libraries and packages in package java and javax.
- By default, a class resides in the *anonymous package*.
- To put it elsewhere, use a package declaration at start of file, as in

```
package database;      or      package ucb.util;
```

- Oracle's javac uses convention that class C in package P1.P2 goes in subdirectory P1/P2 of any other directory in the *class path*.
- Unix example:

```
$ export CLASSPATH=.:$HOME/java-utils:$MASTERDIR/lib/classes/junit.jar  
$ java junit.textui.TestRunner MyTests
```

Searches for TestRunner.class in ./junit/textui, ~/java-utils/junit/textui and finally looks for junit/textui/TestRunner.class in the junit.jar file (which is a single file that is a special compressed archive of an entire directory of files).

Access Modifiers

- Access modifiers (`private`, `public`, `protected`) do not add anything to the power of Java.
- Basically allow a programmer to declare which classes are supposed to need to access ("know about") what declarations.
- In Java, are also part of security—prevent programmers from accessing things that would "break" the runtime system.
- Accessibility always determined by static types.
 - To determine correctness of writing `x.f()`, look at the definition of `f` in the *static type* of `x`.
 - Why the static type? Because the rules are supposed to be enforced by the compiler, which only knows static types of things (static types don't depend on what happens at execution time).

The Access Rules: Public

- Accessibility of a member depends on (1) how the member's declaration is qualified and (2) where it is being accessed.
- `C1`, `C2`, `C3`, and `C4` are distinct classes.
- Class `C2a` is either class `C2` itself or a subtype of `C2`.

```
package P1;
public class C1 ... {
    // M is a method, field, ...
    public int M ...
    void h(C1 x)
        { ... x.M ... } // OK.
}
```

```
package P1;
public class C4 ... {
    void p(C1 x)
        { ... x.M ... } // OK.
}
```

```
package P2;
class C2 extends C3 {
    void f(P1.C1 x) {... x.M ...} // OK
    void g(C2a y) {... y.M ... } // OK
}
```

Public members are available everywhere.

The Access Rules: Private

- `C1`, `C2`, and `C4` are distinct classes.
- Class `C2a` is either class `C2` itself or a subtype of `C2`.

```
package P1;
public class C1 ... {
    // M is a method, field, ...
    private int M ...
    void h(C1 x)
        { ... x.M ... } // OK.
}
```

```
package P1;
public class C4 ... {
    void p(C1 x)
        { ... x.M ... } // ERROR.
}
```

```
package P2;
class C2 extends C1 {
    void f(P1.C1 x) {... x.M ...} // ERROR
    void g(C2a y) {... y.M ... } // ERROR
}
```

Private members are available only within the text of the same class, even for subtypes.

Private to a Class, Not an Object

- Consider a simple list implementation using linking:

```
public class LinkedList {  
    private Node _head;  
    private static class Node {  
        Object _head;  
        Node _tail;  
    }  
  
    /** Swap the contents of this LinkedList and OTHER. */  
    public void swap(LinkedList other) {  
        LinkedList tmp = other._head; other._head = _head; _head = tmp;  
    }...  
}
```

- This is an obvious and correct implementation of swap, *even though* it fetches a private member of a LinkedList object *other than* this.
- It is perfectly OK to do so by the rules, because swap is in the text of LinkedList and therefore has access to all private fields.
- That is, “privacy” does *not* apply to individual objects, only to classes.

The Access Rules: Package Private

- `C1`, `C2`, and `C4` are distinct classes.
- Class `C2a` is either class `C2` itself or a subtype of `C2`.

```
package P1;
public class C1 ... {
    // M is a method, field, ...
    int M ...
    void h(C1 x)
        { ... x.M ... } // OK.
}
```

```
package P1;
public class C4 ... {
    void p(C1 x)
        { ... x.M ... } // OK.
}
```

```
package P2;
class C2 extends C1 {
    void f(P1.C1 x) {... x.M ...} // ERROR
    void g(C2a y) {... y.M ... } // ERROR
}
```

Package Private members are available only within the same package (even for subtypes).

The Access Rules: Protected

- C_1, C_2 , and C_4 are distinct classes.
- Class C_{2a} is either class C_2 itself or a subtype of C_2 .

```
package P1;
public class C1 ... {
    // M is a method, field, ...
    protected int M ...
    void h(C1 x)
        { ... x.M ... } // OK.
}
```

```
package P1;
public class C4 ... {
    void p(C1 x)
        { ... x.M ... } // OK.
}
```

```
package P2;
class C2 extends C1 {
    void f(P1.C1 x) {... x.M ...} // ERROR
        // (x's type is not subtype of C2.)
    void g(C2a y) {... y.M ... } // OK
    void g2()      {... M ... } // OK (this.M)
}
```

Protected members of C_1 are available within P_1 , as for package private. Outside P_1 , they are available within subtypes of C_1 such as C_2 , but only if accessed from expressions whose static types are subtypes of C_2 .

What May be Controlled

- Classes and interfaces that are not nested may be public or package private.
- Members—fields, methods, constructors, and (later) nested types—may have any of the four access levels.
- May *override* a method only with one that has *at least* as permissive an access level. Reason: avoid inconsistency:

```
package P1;
public class C1 {
    public int f() { ... }
}

public class C2 extends C1 {
    // Actually a compiler error; pretend
    // it's not and see what happens
    int f() { ... }
}
```

```
package P2;
class C3 {
    void g(C2 y2) {
        C1 y1 = y2
        y2.f(); // Bad???
        y1.f(); // OK??!?
    }
}
```

That is, there's no point in restricting *C2.f*, because access control depends on static types, and *C1.f* is public.

Intentions of this Design

- `public` declarations represent *specifications*—what clients of a package are supposed to rely on.
- `package private` declarations are part of the *implementation* of a class that must be known to other classes that assist in the implementation.
- `protected` declarations are part of the implementation that subtypes may need, but that clients of the subtypes generally won't.
- `private` declarations are part of the implementation of a class that only that class needs.

Quick Quiz

```
package SomePack;                                // Anonymous package
public class A1 {
    int f1() {
        A1 a = ...
        a.x1 = 3; // OK?
    }
    protected int y1;
    private int x1;
}

class A2 {
    void g(SomePack.A1 x) {
        x.f1(); // OK?
        x.y1 = 3; // OK?
    }
}

class B2 extends SomePack.A1 {
    void h(SomePack.A1 x) {
        x.f1(); // OK?
        x.y1 = 3; // OK?
        f1(); // OK?
        y1 = 3; // OK?
        x1 = 3; // OK?
    }
}
```

- Note: Last three lines of `h` have implicit `this.`'s in front. Static type of `this` is `B2`.

Quick Quiz

```
package SomePack;                                // Anonymous package
public class A1 {
    int f1() {
        A1 a = ...
        a.x1 = 3; // OK
    }
    protected int y1;
    private int x1;
}

class A2 {
    void g(SomePack.A1 x) {
        x.f1(); // OK?
        x.y1 = 3; // OK?
    }
}

class B2 extends SomePack.A1 {
    void h(SomePack.A1 x) {
        x.f1(); // OK?
        x.y1 = 3; // OK?
        f1(); // OK?
        y1 = 3; // OK?
        x1 = 3; // OK?
    }
}
```

- Note: Last three lines of `h` have implicit `this.`'s in front. Static type of `this` is `B2`.

Quick Quiz

```
package SomePack;                                // Anonymous package
public class A1 {
    int f1() {
        A1 a = ...
        a.x1 = 3; // OK
    }
    protected int y1;
    private int x1;
}

class A2 {
    void g(SomePack.A1 x) {
        x.f1(); // ERROR
        x.y1 = 3; // OK?
    }
}

class B2 extends SomePack.A1 {
    void h(SomePack.A1 x) {
        x.f1(); // OK?
        x.y1 = 3; // OK?
        f1(); // OK?
        y1 = 3; // OK?
        x1 = 3; // OK?
    }
}
```

- Note: Last three lines of `h` have implicit `this.`'s in front. Static type of `this` is `B2`.

Quick Quiz

```
package SomePack;                                // Anonymous package
public class A1 {
    int f1() {
        A1 a = ...
        a.x1 = 3; // OK
    }
    protected int y1;
    private int x1;
}

class A2 {
    void g(SomePack.A1 x) {
        x.f1(); // ERROR
        x.y1 = 3; // ERROR
    }
}

class B2 extends SomePack.A1 {
    void h(SomePack.A1 x) {
        x.f1(); // OK?
        x.y1 = 3; // OK?
        f1(); // OK?
        y1 = 3; // OK?
        x1 = 3; // OK?
    }
}
```

- Note: Last three lines of `h` have implicit `this.`'s in front. Static type of `this` is `B2`.

Quick Quiz

```
package SomePack;                                // Anonymous package
public class A1 {
    int f1() {
        A1 a = ...
        a.x1 = 3; // OK
    }
    protected int y1;
    private int x1;
}

class A2 {
    void g(SomePack.A1 x) {
        x.f1(); // ERROR
        x.y1 = 3; // ERROR
    }
}

class B2 extends SomePack.A1 {
    void h(SomePack.A1 x) {
        x.f1(); // ERROR
        x.y1 = 3; // OK?
        f1(); // OK?
        y1 = 3; // OK?
        x1 = 3; // OK?
    }
}
```

- Note: Last three lines of `h` have implicit `this.`'s in front. Static type of `this` is `B2`.

Quick Quiz

```
package SomePack;                                // Anonymous package
public class A1 {
    int f1() {
        A1 a = ...
        a.x1 = 3; // OK
    }
    protected int y1;
    private int x1;
}

class A2 {
    void g(SomePack.A1 x) {
        x.f1(); // ERROR
        x.y1 = 3; // ERROR
    }
}

class B2 extends SomePack.A1 {
    void h(SomePack.A1 x) {
        x.f1(); // ERROR
        x.y1 = 3; // OK?
        f1(); // ERROR
        y1 = 3; // OK?
        x1 = 3; // OK?
    }
}
```

- Note: Last three lines of `h` have implicit `this.`'s in front. Static type of `this` is `B2`.

Quick Quiz

```
package SomePack;                                // Anonymous package
public class A1 {
    int f1() {
        A1 a = ...
        a.x1 = 3; // OK
    }
    protected int y1;
    private int x1;
}

class A2 {
    void g(SomePack.A1 x) {
        x.f1(); // ERROR
        x.y1 = 3; // ERROR
    }
}

class B2 extends SomePack.A1 {
    void h(SomePack.A1 x) {
        x.f1(); // ERROR
        x.y1 = 3; // OK?
        f1(); // ERROR
        y1 = 3; // OK
        x1 = 3; // OK?
    }
}
```

- Note: Last three lines of `h` have implicit `this.`'s in front. Static type of `this` is `B2`.

Quick Quiz

```
package SomePack;                                // Anonymous package
public class A1 {
    int f1() {
        A1 a = ...
        a.x1 = 3; // OK
    }
    protected int y1;
    private int x1;
}

class A2 {
    void g(SomePack.A1 x) {
        x.f1(); // ERROR
        x.y1 = 3; // ERROR
    }
}

class B2 extends SomePack.A1 {
    void h(SomePack.A1 x) {
        x.f1(); // ERROR
        x.y1 = 3; // OK?
        f1(); // ERROR
        y1 = 3; // OK
        x1 = 3; // ERROR
    }
}
```

- Note: Last three lines of `h` have implicit `this.`'s in front. Static type of `this` is `B2`.

Quick Quiz

```
package SomePack;                                // Anonymous package
public class A1 {
    int f1() {
        A1 a = ...
        a.x1 = 3; // OK
    }
    protected int y1;
    private int x1;
}

class A2 {
    void g(SomePack.A1 x) {
        x.f1(); // ERROR
        x.y1 = 3; // ERROR
    }
}

class B2 extends SomePack.A1 {
    void h(SomePack.A1 x) {
        x.f1(); // ERROR
        x.y1 = 3; // ERROR
        f1(); // ERROR
        y1 = 3; // OK
        x1 = 3; // ERROR
    }
}
```

- Note: Last three lines of `h` have implicit `this.`'s in front. Static type of `this` is `B2`.

Access Control Static Only

“Public” and “private” don’t apply to dynamic types; it is possible to call methods in objects of types you can’t name:

```
package utils;
/** A Set of things. */
public interface Collector {
    void add(Object x);
}

-----
package utils;
public class Utils {
    public static Collector concat() {
        return new Concatenator();
    }
}

/** NON-PUBLIC class that collects strings. */
class Concatenator implements Collector {
    StringBuffer stuff = new StringBuffer();
    int n = 0;
    public void add(Object x) { stuff.append(x); n += 1; }
    public Object value() { return stuff.toString(); }
}

| package mystuff;
|
| class User {
|     utils.Collector c =
|         utils.Utils.concat();
|
|     c.add("foo"); // OK
|     ... c.value(); // ERROR
|     ((utils.Concatenator) c).value()
|                     // ERROR
|-----
```

CS61B Lecture #15: Integers

Announcements:

- Project 1 checkpoint due Friday.
- Please use gitbug (see the Gitbugs tab on the website) to submit requests for help debugging projects, homeworks, etc. This can be a great deal more efficient than office hours or Piazza. In particular, it helps to make sure we have all the information needed to help you.

Integer Types and Literals

Type	Bits	Signed?	Literals
byte	8	Yes	Cast from int: (byte) 3
short	16	Yes	None. Cast from int: (short) 4096
char	16	No	'a' // (char) 97 '\n' // newline ((char) 10) '\t' // tab ((char) 8) '\\' // backslash 'A', '\101', '\u0041' // == (char) 65
int	32	Yes	123 0100 // Octal for 64 0x3f, 0xffffffff // Hexadecimal 63, -1 (!)
long	64	Yes	123L, 01000L, 0x3fL 1234567891011L

- Negative numerals are just negated (positive) literals.
- “ N bits” means that there are 2^N integers in the domain of the type:
 - If signed, range of values is $-2^{N-1} \dots 2^{N-1} - 1$.
 - If unsigned, only non-negative numbers, and range is $0 \dots 2^N - 1$.

Overflow

- **Problem:** How do we handle overflow, such as occurs in $10000 * 10000 * 10000$?
- Some languages throw an exception (Ada), some give undefined results (C, C++)
- Java **defines** the result of any arithmetic operation or conversion on integer types to “wrap around”—**modular arithmetic**.
- That is, the “next number” after the largest in an integer type is the smallest (like “clock arithmetic”).
- E.g., (byte) 128 == (byte) (127+1) == (byte) -128
- In general,
 - If the result of some arithmetic subexpression is supposed to have type T , an n -bit integer type,
 - then we compute the real (mathematical) value, x ,
 - and yield a number, x' , that is in the range of T , and that is equivalent to x modulo 2^n .
 - (That means that $x - x'$ is a multiple of 2^n .)

Modular Arithmetic

- Define $a \equiv b \pmod{n}$ to mean that $a - b = kn$ for some integer k .
- Define the binary operation $a \bmod n$ as the value b such that $a \equiv b \pmod{n}$ and $0 \leq b < n$ for $n > 0$. (Can be extended to $n \leq 0$ as well, but we won't bother with that here.) This is *not* the same as Java's `%` operation.
- Various facts: (Here, let a' denote $a \bmod n$).

$$\begin{aligned}a'' &= a' \\(a + b)' &= (a' + b')' = (a + b)' \\(a - b)' &= (a' + (-b)')' = (a' - b')' \\(a \cdot b)' &= (a' \cdot b')' \\(a^k)' &= ((a')^k)' = (a' \cdot (a^{k-1})')', \text{ for } k > 0.\end{aligned}$$

Modular Arithmetic: Examples

- (byte) $(64*8)$ yields 0, since $512 - 0 = 2 \times 2^8$.
- (byte) $(64*2)$ and (byte) $(127+1)$ yield -128, since $128 - (-128) = 1 \times 2^8$.
- (byte) $(101*99)$ yields 15, since $9999 - 15 = 39 \times 2^8$.
- (byte) $(-30*13)$ yields 122, since $-390 - 122 = -2 \times 2^8$.
- (char) (-1) yields $2^{16} - 1$, since $-1 - (2^{16} - 1) = -1 \times 2^{16}$.

Modular Arithmetic and Bits

- Why wrap around?
- Java's definition is the natural one for a machine that uses binary arithmetic.
- For example, consider bytes (8 bits):

Decimal	Binary
101	1100101
$\times 99$	1100011
9999	100111 00001111
– 9984	100111 00000000
15	00001111

- In general, bit n , counting from 0 at the right, corresponds to 2^n .
- The bits to the left of the vertical bars therefore represent multiples of $2^8 = 256$.
- So throwing them away is the same as arithmetic modulo 256.

Negative numbers

- Why this representation for -1?

$$\begin{array}{r} 1 \quad 00000001_2 \\ + -1 \quad 11111111_2 \\ = 0 \quad | 00000000_2 \end{array}$$

Only 8 bits in a byte, so bit 8 falls off, leaving 0.

- The truncated bit is in the 2^8 place, so throwing it away gives an equal number modulo 2^8 . All bits to the left of it are also divisible by 2^8 .
- On unsigned types (`char`), arithmetic is the same, but we choose to represent only non-negative numbers modulo 2^{16} :

$$\begin{array}{r} 1 \quad 0000000000000001_2 \\ + 2^{16} - 1 \quad 1111111111111111_2 \\ = 2^{16} + 0 \quad | 0000000000000000_2 \end{array}$$

Conversion

- In general Java will silently convert from one type to another if this makes sense and no information is lost from value.
- Otherwise, cast explicitly, as in (byte) x.
- Hence, given

```
byte aByte; char aChar; short aShort; int anInt; long aLong;

// OK:
aShort = aByte; anInt = aByte; anInt = aShort;
anInt = aChar; aLong = anInt;

// Not OK, might lose information:
anInt = aLong; aByte = anInt; aChar = anInt; aShort = anInt;
aShort = aChar; aChar = aShort; aChar = aByte;

// OK by special dispensation:
aByte = 13;      // 13 is compile-time constant
aByte = 12+100 // 112 is compile-time constant
```

Promotion

- Arithmetic operations (+, *, ...) *promote* operands as needed.
- Promotion is just implicit conversion.
- For integer operations,
 - if any operand is **long**, promote both to **long**.
 - otherwise promote both to **int**.
- So,

```
aByte + 3 == (int) aByte + 3    // Type int
aLong + 3 == aLong + (long) 3   // Type long
'A' + 2 == (int) 'A' + 2       // Type int
aByte = aByte + 1              // ILLEGAL (why?)
```

- But fortunately,

```
aByte += 1;                  // Defined as aByte = (byte) (aByte+1)
```

- Common example:

```
// Assume aChar is an upper-case letter
char lowerCaseChar = (char) ('a' + aChar - 'A'); // why cast?
```

Bit twiddling

- Java (and C, C++) allow for handling integer types as sequences of bits. No “conversion to bits” needed: they already are.
- Operations and their uses:

Mask	Set	Flip	Flip all
00101100 & 10100111	00101100 10100111	00101100 ^ 10100111	~ 10100111
	00100100 10101111	10001011	01011000

- Shifting:

Left	Arithmetic Right	Logical Right
10101101 << 3	10101101 >> 3	10101100 >>> 3
01101000	11110101	00010101
	01010110 >> 3	
	00001010	

(-1) >>> 29?

- What is:
 $x << n?$
 $x >> n?$
 $(x >>> 3) \& ((1<<5)-1)?$

Bit twiddling

- Java (and C, C++) allow for handling integer types as sequences of bits. No “conversion to bits” needed: they already are.
- Operations and their uses:

Mask	Set	Flip	Flip all
00101100 & 10100111	00101100 10100111	00101100 ^ 10100111	~ 10100111
	00100100 10101111	10001011	01011000

- Shifting:

Left	Arithmetic Right	Logical Right
10101101 << 3	10101101 >> 3	10101100 >>> 3
01101000	11110101	00010101
	01010110 >> 3	
	00001010	
(-1) >>> 29?		= 7.
• What is: $x << n?$ $x >> n?$ $(x >>> 3) \& ((1<<5)-1)?$		

Bit twiddling

- Java (and C, C++) allow for handling integer types as sequences of bits. No “conversion to bits” needed: they already are.
- Operations and their uses:

Mask	Set	Flip	Flip all
00101100	00101100	00101100	
& 10100111	10100111	^ 10100111	~ 10100111
00100100	10101111	10001011	01011000

- Shifting:

Left	Arithmetic Right	Logical Right
10101101 << 3	10101101 >> 3	10101100 >>> 3
01101000	11110101	00010101
	01010110 >> 3	
	00001010	

(-1) >>> 29?	= 7.
$x << n?$	
$x >> n?$	
$(x >>> 3) \& ((1<<5)-1)?$	$= x \cdot 2^n.$

Bit twiddling

- Java (and C, C++) allow for handling integer types as sequences of bits. No “conversion to bits” needed: they already are.
- Operations and their uses:

Mask	Set	Flip	Flip all
00101100	00101100	00101100	
& 10100111	10100111	^ 10100111	~ 10100111
00100100	10101111	10001011	01011000

- Shifting:

Left	Arithmetic Right	Logical Right
10101101 << 3	10101101 >> 3	10101100 >>> 3
01101000	11110101	00010101
	01010110 >> 3	
	00001010	

- What is: $x \ll n?$
 - $x \gg n?$
 - $(x \ggg 3) \& ((1\ll 5)-1)?$
- $= 7.$
 $= x \cdot 2^n.$
 $= \lfloor x/2^n \rfloor$ (i.e., rounded down).

Bit twiddling

- Java (and C, C++) allow for handling integer types as sequences of bits. No “conversion to bits” needed: they already are.
- Operations and their uses:

Mask	Set	Flip	Flip all
00101100 & 10100111	00101100 10100111	00101100 ^ 10100111	~ 10100111
	00100100 10101111	10001011	01011000

- Shifting:

Left	Arithmetic Right	Logical Right
10101101 << 3	10101101 >> 3	10101100 >>> 3
01101000	11110101	00010101
	01010110 >> 3	
	00001010	

- What is:
 $x << n?$
 $x >> n?$
 $(x >>> 3) \& ((1<<5)-1)?$
- $= 7.$
 $= x \cdot 2^n.$
 $= \lfloor x/2^n \rfloor$ (i.e., rounded down).
5-bit integer, bits 3-7 of x .

CS61B Lecture #16: Complexity

What Are the Questions?

- Cost is a principal concern throughout engineering:
"An engineer is someone who can do for a dime what any fool can do for a dollar."
- Cost can mean
 - *Operational cost* (for programs, time to run, space requirements).
 - *Development costs*: How much engineering time? When delivered?
 - *Maintenance costs*: Upgrades, bug fixes.
 - *Costs of failure*: How robust? How safe?
- Is this program *fast enough*? Depends on:
 - *For what purpose*;
 - *For what input data*.
- How much *space* (memory, disk space)?
 - Again depends on what input data.
- How will it *scale*, as input gets big?

Enlightening Example

Problem: Scan a text corpus (say 10^9 bytes or so), and find and print the 20 most frequently used words, together with counts of how often they occur.

- Solution 1 (Knuth): Heavy-Duty data structures
 - Hash Trie implementation, randomized placement, pointers galore, several pages long.
- Solution 2 (Doug McIlroy): UNIX shell script:

```
tr -c -s '[:alpha:]' '[\n*]' < FILE | \
sort | \
uniq -c | \
sort -n -r -k 1,1 | \
sed 20q
```

- Which is better?
 - #1 is much faster,
 - but #2 took 5 minutes to write and processes 1GB in ≈ 256 sec.
 - I pick #2.
- In very many cases, almost anything will do: **Keep It Simple.**

Cost Measures (Time)

- *Wall-clock or execution time*

- You can do this at home:

```
time java FindPrimes 1000
```

- Advantages: easy to measure, meaning is obvious.
 - Appropriate where time is critical (real-time systems, e.g.).
 - Disadvantages: applies only to specific data set, compiler, machine, etc.

- *Dynamic statement counts* of # of times statements are executed:

- Advantages: more general (not sensitive to speed of machine).
 - Disadvantages: doesn't tell you actual time, still applies only to specific data sets.

- *Symbolic execution times*:

- That is, *formulas* for execution times as functions of input size.
 - Advantages: applies to all inputs, makes scaling clear.
 - Disadvantage: practical formula must be approximate, may tell very little about actual time.

Asymptotic Cost

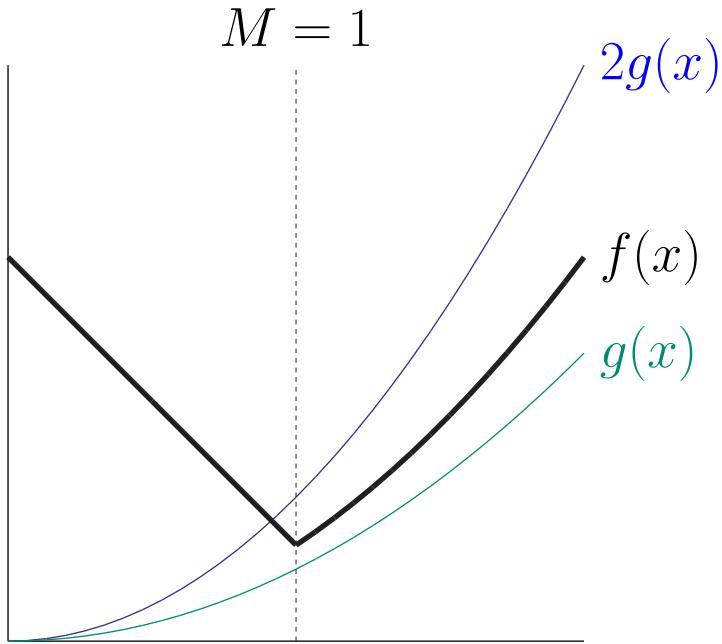
- Symbolic execution time lets us see *shape* of the cost function.
- Since we are approximating anyway, pointless to be precise about certain things:
 - *Behavior on small inputs*:
 - * Can always pre-calculate some results.
 - * Times for small inputs not usually important.
 - * Often more interested in *asymptotic behavior* as input size becomes very large.
 - *Constant factors* (as in "off by factor of 2"):
 - * Just changing machines causes constant-factor change.
- How to abstract away from (i.e., ignore) these things?

Handy Tool: Order Notation

- Idea: Don't try to produce specific functions that specify size, but rather *families of functions with similarly behaved magnitudes*.
- Then say something like " f is bounded by g if it is in g 's family."
- For any function $g(x)$, the functions $2g(x)$, $0.5g(x)$, or for any $K > 0$, $K \cdot g(x)$, all have the same "shape". So put all of them into g 's family.
- Any function $h(x)$ such that $h(x) = K \cdot g(x)$ for $x > M$ (for some constant M) has g 's shape "except for small values." So put all of these in g 's family.
- For upper limits, throw in all functions whose absolute value is everywhere \leq *some member of* g 's family. Call this set $O(g)$ or $O(g(n))$.
- Or, for lower limits, throw in all functions whose absolute value is everywhere \geq some member of g 's family. Call this set $\Omega(g)$.
- Finally, define $\Theta(g) = O(g) \cap \Omega(g)$ —the set of functions *bracketed in magnitude* by two members of g 's family.

Big Oh

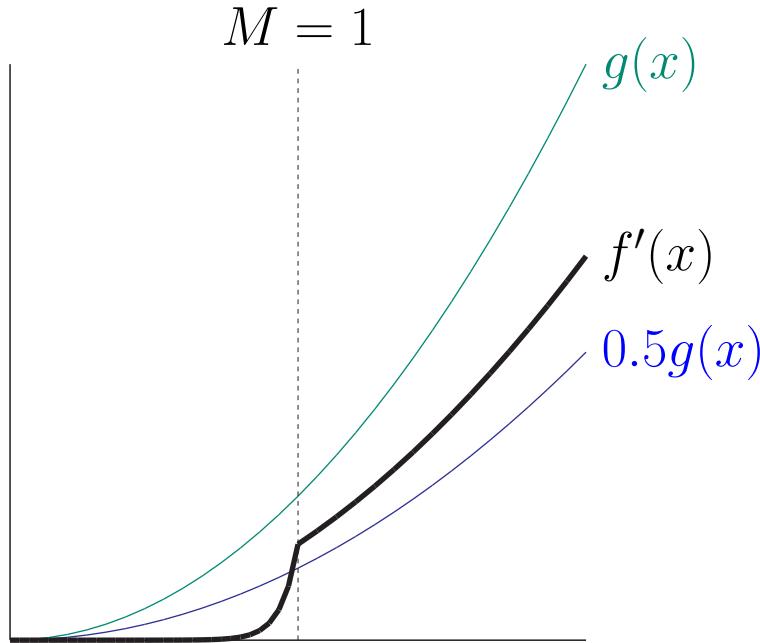
- Goal: Specify bounding from above.



- Here, $f(x) \leq 2g(x)$ as long as $x > 1$,
- So $f(x)$ is in g 's "bounded-above family," written
$$f(x) \in O(g(x)),$$
- ... even though (in this case) $f(x) > g(x)$ everywhere.

Big Omega

- Goal: Specify bounding from below:



- Here, $f'(x) \geq \frac{1}{2}g(x)$ as long as $x > 1$,
- So $f'(x)$ is in g 's "bounded-below family," written

$$f'(x) \in \Omega(g(x)),$$

- ... even though $f(x) < g(x)$ everywhere.

Big Theta

- In the two previous slides, we not only have $f(x) \in O(g(x))$ and $f'(x) \in \Omega(g(x)), \dots$
- ... but also $f(x) \in \Omega(g(x))$ and $f'(x) \in O(g(x)).$
- We can summarize this all by saying $f(x) \in \Theta(g(x))$ and $f'(x) \in \Theta(g(x)).$

Aside: Various Mathematical Pedantry

- Technically, if I am going to talk about $O(\cdot)$, $\Omega(\cdot)$ and $\Theta(\cdot)$ as sets of functions, I really should write, for example,

$$f \in O(g) \quad \text{instead of} \quad f(x) \in O(g(x))$$

- In effect, $f(x) \in O(g(x))$ is short for $\lambda x. f(x) \in O(\lambda x. g(x))$.
- The standard notation outside this course, in fact, is $f(x) = O(g(x))$, but personally, I think that's a serious abuse of notation.

How We Use Order Notation

- Elsewhere in mathematics, you'll see $O(\dots)$, etc., used generally to specify bounds on functions.
- For example,

$$\pi(N) = \Theta\left(\frac{N}{\ln N}\right)$$

which I would prefer to write

$$\pi(N) \in \Theta\left(\frac{N}{\ln N}\right)$$

(Here, $\pi(N)$ is the number of primes less than or equal to N .)

- Also, you'll see things like

$$f(x) = x^3 + x^2 + O(x) \quad (\text{or } f(x) \in x^3 + x^2 + O(x)),$$

meaning that $f(x) = x^3 + x^2 + g(x)$ where $g(x) \in O(x)$.

- For our purposes, the functions we will be bounding will be **cost functions**: functions that measure the amount of execution time or the amount of space required by a program or algorithm.

Why It Matters

- Computer scientists often talk as if constant factors didn't matter at all, only the difference of $\Theta(N)$ vs. $\Theta(N^2)$.
- In reality they do matter, but at some point, constants always get swamped.

n	$16 \lg n$	\sqrt{n}	n	$n \lg n$	n^2	n^3	2^n
2	16	1.4	2	2	4	8	4
4	32	2	4	8	16	64	16
8	48	2.8	8	24	64	512	256
16	64	4	16	64	256	4,096	65,636
32	80	5.7	32	160	1024	32,768	4.2×10^9
64	96	8	64	384	4,096	262,144	1.8×10^{19}
128	112	11	128	896	16,384	2.1×10^9	3.4×10^{38}
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
1,024	160	32	1,024	10,240	1.0×10^6	1.1×10^9	1.8×10^{308}
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
2^{20}	320	1024	1.0×10^6	2.1×10^7	1.1×10^{12}	1.2×10^{18}	$6.7 \times 10^{315,652}$

- For example: replace column n^2 with $10^6 \cdot n^2$ and it still becomes dominated by 2^n .

Some Intuition on Meaning of Growth

- How big a problem can you solve in a given time?
- In the following table, left column shows time in microseconds to solve a given problem as a function of problem size N .
- Entries show the *size of problem* that can be solved in a second, hour, month (31 days), and century, for various relationships between time required and problem size.
- N = problem size.

Time (μ sec) for problem size N	1 second	Max N Possible in 1 hour	1 month	1 century
$\lg N$	10^{300000}	$10^{10000000000}$	$10^8 \cdot 10^{11}$	$10^{10} \cdot 10^{14}$
N	10^6	$3.6 \cdot 10^9$	$2.7 \cdot 10^{12}$	$3.2 \cdot 10^{15}$
$N \lg N$	63000	$1.3 \cdot 10^8$	$7.4 \cdot 10^{10}$	$6.9 \cdot 10^{13}$
N^2	1000	60000	$1.6 \cdot 10^6$	$5.6 \cdot 10^7$
N^3	100	1500	14000	150000
2^N	20	32	41	51

Using the Notation

- Can use this order notation for any kind of real-valued function.
- We will use them to describe cost functions. Example:

```
/** Find position of X in list L, or -1 if not found. */
int find(List L, Object X) {
    int c;
    for (c = 0; L != null; L = L.next, c += 1)
        if (X.equals(L.head)) return c;
    return -1;
}
```

- Choose representative operation: number of `.equals` tests.
- If N is length of L , then loop does *at most* N tests: *worst-case time* is N tests.
- In fact, total # of instructions executed is roughly proportional to N in the worst case, so can also say worst-case time is $O(N)$, regardless of units used to measure.
- Use $N > M$ provision (in defn. of $O(\cdot)$) to ignore empty list.

Be Careful

- It's also true that the worst-case time is $O(N^2)$, since $N \in O(N^2)$ also: Big-Oh bounds are loose.
- The worst-case time is $\Omega(N)$, since $N \in \Omega(N)$, but that does *not* mean that the loop *always* takes time N , or even $K \cdot N$ for some K .
- Instead, we are just saying something about the *function* that maps N into the *largest possible* time required to process any array of length N .
- To say as much as possible about our worst-case time, we should try to give a Θ bound: in this case, we can: $\Theta(N)$.
- But again, that still tells us nothing about *best-case* time, which happens when we find X at the beginning of the loop. Best-case time is $\Theta(1)$.

Effect of Nested Loops

- Nested loops often lead to polynomial bounds:

```
for (int i = 0; i < A.length; i += 1)
    for (int j = 0; j < A.length; j += 1)
        if (i != j && A[i] == A[j])
            return true;
return false;
```

- Clearly, time is $O(N^2)$, where $N = A.length$. *Worst-case time* is $\Theta(N^2)$.

Constant Factor Speed-Up

- Previous loop is inefficient. This one is considerably faster:

```
for (int i = 0; i < A.length; i += 1)
    for (int j = i+1; j < A.length; j += 1)
        if (A[i] == A[j]) return true;
return false;
```

- Now worst-case time is proportional to

$$N - 1 + N - 2 + \dots + 1 = N(N - 1)/2,$$

which is a constant-factor improvement.

- But still,

$$N(N - 1)/2 \in \Theta(N^2).$$

so the asymptotic time is unchanged by the constant-factor speed-up.

Recursion and Recurrences: Fast Growth

- Silly example of recursion. In the worst case, both recursive calls happen:

```
/** True iff X is a substring of S */
boolean occurs(String S, String X) {
    if (S.equals(X)) return true;
    if (S.length() <= X.length()) return false;
    return
        occurs(S.substring(1), X) ||
        occurs(S.substring(0, S.length()-1), X);
}
```

- Define $C(N)$ to be the worst-case cost of `occurs(S,X)` for S of length N , X of fixed size N_0 , measured in # of calls to `occurs`. Then

$$C(N) = \begin{cases} 1, & \text{if } N \leq N_0, \\ 2C(N-1) + 1 & \text{if } N > N_0 \end{cases}$$

- So $C(N)$ grows exponentially:

$$\begin{aligned} C(N) &= 2C(N-1) + 1 = 2(2C(N-2) + 1) + 1 = \dots = \underbrace{2(\dots 2 \cdot 1 + 1)}_{N-N_0} + \dots + 1 \\ &= 2^{N-N_0} + 2^{N-N_0-1} + 2^{N-N_0-2} + \dots + 1 = 2^{N-N_0+1} - 1 \in \Theta(2^N) \end{aligned}$$

Binary Search: Slow Growth

```
/** True X iff is an element of S[L .. U]. Assumes
 *  S in ascending order, 0 <= L <= U-1 < S.length. */
boolean isIn(String X, String[] S, int L, int U) {
    if (L > U) return false;
    int M = (L+U)/2;
    int direct = X.compareTo(S[M]);
    if (direct < 0) return isIn(X, S, L, M-1);
    else if (direct > 0) return isIn(X, S, M+1, U);
    else return true;
}
```

- Here, worst-case time, $C(D)$, (as measured by # of calls to `.compareTo`), depends on size $D = U - L + 1$.
- We eliminate $S[M]$ from consideration each time and look at half the rest. Assume $D = 2^k - 1$ for simplicity, so:

$$\begin{aligned} C(D) &= \begin{cases} 0, & \text{if } D \leq 0, \\ 1 + C((D-1)/2), & \text{if } D > 0. \end{cases} \\ &= \underbrace{1 + 1 + \dots + 1}_{k} + 0 \\ &= k = \lg(D+1) \in \Theta(\lg D) \end{aligned}$$

Another Typical Pattern: Merge Sort

```
List sort(List L) {  
    if (L.length() < 2) return L;  
    Split L into L0 and L1 of about equal size;  
    L0 = sort(L0); L1 = sort(L1);  
    return Merge of L0 and L1  
}
```

Merge ("combine into a single ordered list") takes time proportional to size of its result.

- Assuming that size of L is $N = 2^k$, worst-case cost function, $C(N)$, counting just merge time (which is proportional to # items merged):

$$\begin{aligned} C(N) &= \begin{cases} 0, & \text{if } N < 2; \\ 2C(N/2) + N, & \text{if } N \geq 2. \end{cases} \\ &= 2(2C(N/4) + N/2) + N \\ &= 4C(N/4) + N + N \\ &= 8C(N/8) + N + N + N \\ &= N \cdot 0 + \underbrace{N + N + \dots + N}_{k=\lg N} \\ &= N \lg N \end{aligned}$$

- In general, can say it's $\Theta(N \lg N)$ for arbitrary N (not just 2^k).

CS61B Lecture #17

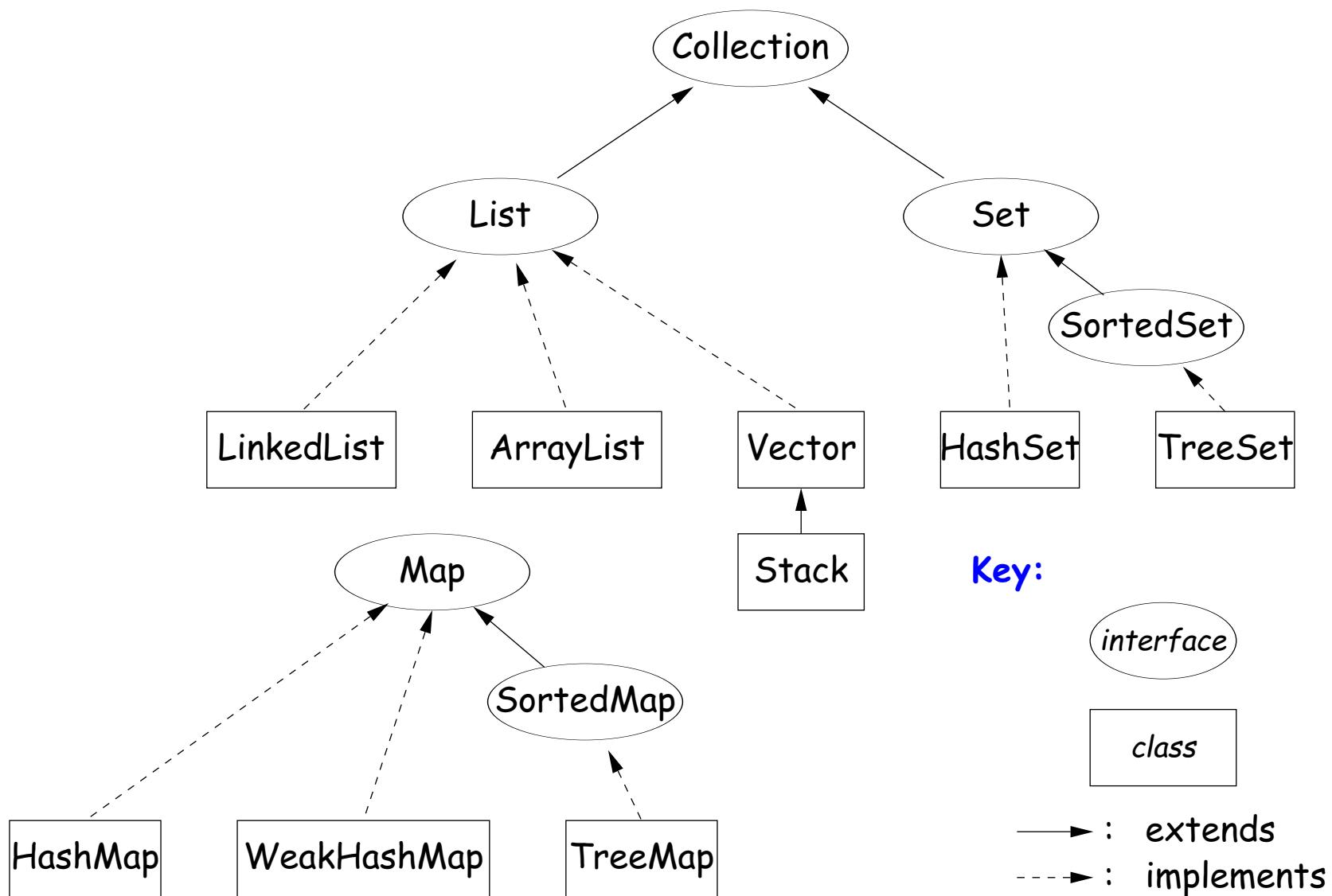
Topics

- Overview of standard Java Collections classes.
 - Iterators, ListIterators
 - Containers and maps in the abstract
- Amortized analysis of implementing lists with arrays.

Data Types in the Abstract

- Most of the time, should *not* worry about implementation of data structures, search, etc.
- What they do for us—their specification—is important.
- Java has several standard types (in `java.util`) to represent collections of objects
 - Six interfaces:
 - * Collection: General collections of items.
 - * List: Indexed sequences with duplication
 - * Set, SortedSet: Collections without duplication
 - * Map, SortedMap: Dictionaries ($\text{key} \mapsto \text{value}$)
 - Concrete classes that provide actual instances: `LinkedList`, `ArrayList`, `HashSet`, `TreeSet`.
 - To make change easier, purists would use the concrete types only for new, interfaces for parameter types, local variables.

Collection Structures in java.util



The Collection Interface

- Collection interface. Main functions promised:
 - Membership tests: contains (\in), containsAll (\subseteq)
 - Other queries: size, isEmpty
 - Retrieval: iterator, toArray
 - *Optional* modifiers: add, addAll, clear, remove, removeAll (set difference), retainAll (intersect)

Side Trip about Library Design: Optional Operations

- Not all Collections need to be modifiable; often makes sense just to get things from them.
- So some operations are optional (add, addAll, clear, remove, removeAll, retainAll)
- The library developers decided to have *all* Collections implement these, but allowed implementations to throw an exception:

UnsupportedOperationException

- An alternative design would have created separate interfaces:

```
interface Collection { contains, containsAll, size, iterator, ... }  
interface Expandable extends Collection { add, addAll }  
interface Shrinkable extends Collection { remove, removeAll, ... }  
interface ModifiableCollection  
    extends Collection, Expandable, Shrinkable { }
```

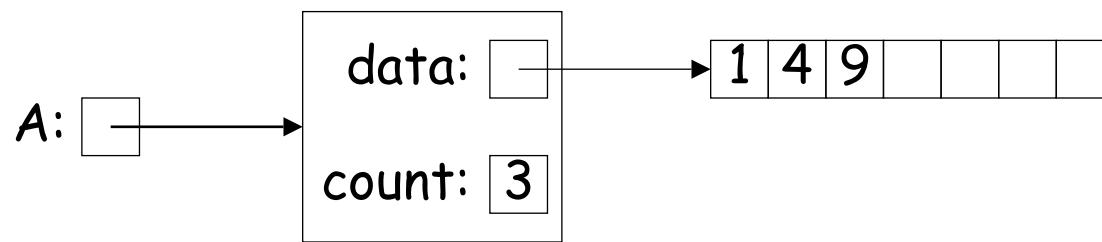
- You'd soon have lots of interfaces. Perhaps that's why they didn't do it that way.

The List Interface

- Extends Collection
- Intended to represent *indexed sequences* (generalized arrays)
- Adds new methods to those of Collection:
 - Membership tests: `indexOf`, `lastIndexOf`.
 - Retrieval: `get(i)`, `listIterator()`, `sublist(B, E)`.
 - Modifiers: add and `addAll` with additional index to say *where* to add. Likewise for removal operations. set operation to go with get.
- Type `ListIterator<Item>` extends `Iterator<Item>`:
 - Adds `previous` and `hasPrevious`.
 - `add`, `remove`, and `set` allow one to iterate through a list, inserting, removing, or changing as you go.
 - **Important Question:** What advantage is there to saying List L rather than `LinkedList L` or `ArrayList L`?

Implementing Lists (I): ArrayLists

- The main concrete types in Java library for interface List are ArrayList and LinkedList:
- As you might expect, an ArrayList, A, uses an array to hold data. For example, a list containing the three items 1, 4, and 9 might be represented like this:



- After adding four more items to A, its data array will be full, and the value of data will have to be replaced with a pointer to a new, bigger array that starts with a copy of its previous values.
- Question: For best performance, how big should this new array be?
- If we increase the size by 1 each time it gets full (or by any constant value), the cost of N additions will scale as $\Theta(N^2)$, which makes ArrayList look much worse than LinkedList (which uses an IntList-like implementation.)

Expanding Vectors Efficiently

- When using array for expanding sequence, best to *double* the size of array to grow it. Here's why.
- If array is size s , doubling its size and moving s elements to the new array takes time proportional to $2s$.
- In all cases, there is an additional $\Theta(1)$ cost for each addition to account for actually assigning the new value into the array.
- When you add up these costs for inserting a sequence of N items, the *total* cost turns out to be proportional to N , as if each addition took constant time, even though some of the additions actually take time proportional to N all by themselves!

Amortized Time

- Suppose that the actual costs of a sequence of N operations are c_0, c_1, \dots, c_{N-1} , which may differ from each other by arbitrary amounts and where $c_i \in O(f(i))$.
- Consider another sequence a_0, a_1, \dots, a_{N-1} , where $a_i \in O(g(i))$.
- If

$$\sum_{0 \leq i < k} a_i \geq \sum_{0 \leq i < k} c_i \text{ for all } k,$$

we say that the operations all run in $O(g(i))$ amortized time.

- That is, the actual cost of a given operation, c_i , may be arbitrarily larger than the amortized time, a_i , as long as the **total** amortized time is always greater than or equal to the total actual time, no matter where the sequence of operations stops—i.e., no matter what k is.
- In cases of interest, the amortized time bounds are much less than the actual individual time bounds: $g(i) \ll f(i)$.
- E.g., for the case of insertion with array doubling, $f(i) \in O(N)$ and $g(i) \in O(1)$.

Amortization: Expanding Vectors (II)

To Insert Item #	Resizing Cost	Cumulative Cost	Resizing Cost per Item	Array Size After Insertions
0	0	0	0	1
1	2	2	1	2
2	4	6	2	4
3	0	6	1.5	4
4	8	14	2.8	8
5	0	14	2.33	8
:	:	:	:	:
7	0	14	1.75	8
8	16	30	3.33	16
:	:	:	:	:
15	0	30	1.88	16
:	:	:	:	:
$2^m + 1$ to $2^{m+1} - 1$	0	$2^{m+2} - 2$	≈ 2	2^{m+1}
2^{m+1}	2^{m+2}	$2^{m+3} - 2$	≈ 4	2^{m+2}

- If we spread out (**amortize**) the cost of resizing, we average at most about 4 time units for resizing on each item: "amortized resizing time is 4 units." Time to add N elements is $\Theta(N)$, **not** $\Theta(N^2)$.

Demonstrating Amortized Time: Potential Method

- To formalize the argument, associate a *potential*, $\Phi_i \geq 0$, to the i^{th} operation that keeps track of “saved up” time from cheap operations that we can “spend” on later expensive ones. Start with $\Phi_0 = 0$.
- Now we pretend that the cost of the i^{th} operation is actually a_i , the *amortized cost*, defined

$$a_i = c_i + \Phi_{i+1} - \Phi_i,$$

where c_i is the real cost of the operation. Or, looking at potential:

$$\Phi_{i+1} = \Phi_i + (a_i - c_i)$$

- On cheap operations, we artificially set $a_i > c_i$ so that we can increase Φ ($\Phi_{i+1} > \Phi_i$).
- On expensive ones, we typically have $a_i \ll c_i$ and greatly decrease Φ (but don’t let it go negative—may not be “overdrawn”).
- We try to do all this so that a_i remains as we desired (e.g., $O(1)$ for expanding array), without allowing $\Phi_i < 0$.
- Requires that we choose a_i so that Φ_i always stays ahead of c_i .

Application to Expanding Arrays

- When adding to our array, the cost, c_i , of adding element # i when the array already has space for it is 1 unit.
- The array does not initially have space when adding items 1, 2, 4, 8, 16, ... —in other words at item 2^n for all $n \geq 0$. So,
 - $c_i = 1$ if $i \geq 0$ and is not a power of 2; and
 - $c_i = 2i + 1$ when i is a power of 2 (copy i items, clear another i items, and then add item # i).
- So on each operation # 2^n we're going to need to have saved up at least $2 \cdot 2^n = 2^{n+1}$ units of potential to cover the expense of expanding the array, and we have this operation and the preceding $2^{n-1} - 1$ operations in which to save up this much potential (everything since the preceding doubling operation).
- So choose $a_0 = 1$ and $a_i = 5$ for $i > 0$. Apply $\Phi_{i+1} = \Phi_i + (a_i - c_i)$, and here is what happens:

i	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
c_i	1	3	5	1	9	1	1	1	17	1	1	1	1	1	1	1	33	1
a_i	1	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5
Φ_i	0	0	2	2	6	2	6	10	14	2	6	10	14	18	22	26	30	2

Pretending each cost is 5 never underestimates true cumulative time.

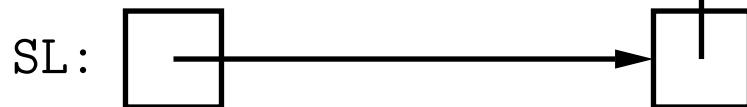
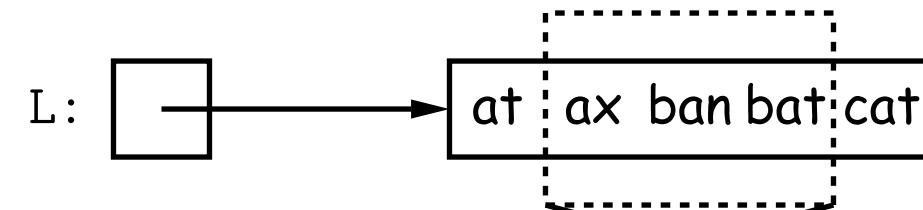
CS61B Lecture #18: Assorted Topics

- Views
- Maps
- More partial implementations
- Array vs. linked: tradeoffs
- Sentinels
- Specialized sequences: stacks, queues, deques
- Circular buffering
- Recursion and stacks
- Adapters

Views

New Concept: A *view* is an alternative presentation of (interface to) an existing object.

- For example, the sublist method is supposed to yield a “view of” part of an existing list:



```
List<String> L = new ArrayList<String>();  
L.add("at"); L.add("ax"); ...  
List<String> SL = L.sublist(1,4);
```

- Example: after `L.set(2, "bag")`, value of `SL.get(1)` is "bag", and after `SL.set(1,"bad")`, value of `L.get(2)` is "bad".
- Example: after `SL.clear()`, `L` will contain only "at" and "cat".
- Small challenge: “How do they do that?!”

Maps

- A Map is a kind of “modifiable function:”

```
package java.util;
public interface Map<Key,Value> {
    Value get(Object key);           // Value at KEY.
    Object put(Key key, Value value); // Set get(KEY) -> VALUE
    ...
}
```

```
Map<String,String> f = new TreeMap<String,String>();
f.put("Paul", "George"); f.put("George", "Martin");
f.put("Dana", "John");
// Now f.get("Paul").equals("George")
//      f.get("Dana").equals("John")
//      f.get("Tom") == null
```

Map Views

```
public interface Map<Key,Value> { // Continuation  
  
    /* Views of Maps */  
  
    /** The set of all keys. */  
    Set<Key> keySet();  
  
    /** The multiset of all values that can be returned by get.  
     * (A multiset is a collection that may have duplicates). */  
    Collection<Value> values();  
  
    /** The set of all(key, value) pairs */  
    Set<Map.Entry<Key,Value>> entrySet();  
}
```

View Examples

Using example from a previous slide:

```
Map<String, String> f = new TreeMap<String, String>();  
f.put("Paul", "George"); f.put("George", "Martin");  
f.put("Dana", "John");
```

we can take various views of f:

```
for (Iterator<String> i = f.keySet().iterator(); i.hasNext();) i.next() ==> Dana, George, Paul
```

// or, more succinctly:

```
for (String name : f.keySet()) name ==> Dana, George, Paul
```

```
for (String parent : f.values()) parent ==> John, Martin, George
```

```
for (Map.Entry<String, String> pair : f.entrySet()) pair ==> (Dana, John), (George, Martin), (Paul, George)
```

```
f.keySet().remove("Dana"); // Now f.get("Dana") == null
```

Simple Banking I: Accounts

Problem: Want a simple banking system. Can look up accounts by name or number, deposit or withdraw, print.

Account Structure

```
class Account {  
    Account(String name, String number, int init) {  
        this.name = name; this.number = number;  
        this.balance = init;  
    }  
    /** Account-holder's name */  
    final String name;  
    /** Account number */  
    final String number;  
    /** Current balance */  
    int balance;  
  
    /** Print THIS on STR in some useful format. */  
    void print(PrintStream str) { ... }  
}
```

Simple Banking II: Banks

```
class Bank {  
    /* These variables maintain mappings of String -> Account. They keep  
     * the set of keys (Strings) in "compareTo" order, and the set of  
     * values (Accounts) is ordered according to the corresponding keys. */  
    SortedMap<String,Account> accounts = new TreeMap<String,Account>();  
    SortedMap<String,Account> names = new TreeMap<String,Account>();  
  
    void openAccount(String name, int initBalance) {  
        Account acc =  
            new Account(name, chooseNumber(), initBalance);  
        accounts.put(acc.number, acc);  
        names.put(name, acc);  
    }  
  
    void deposit(String number, int amount) {  
        Account acc = accounts.get(number);  
        if (acc == null) ERROR(...);  
        acc.balance += amount;  
    }  
    // Likewise for withdraw.
```

Banks (continued): Iterating

Printing out Account Data

```
/** Print out all accounts sorted by number on STR. */
void printByAccount(PrintStream str) {
    // accounts.values() is the set of mapped-to values.  Its
    // iterator produces elements in order of the corresponding keys.
    for (Account account : accounts.values())
        account.print(str);
}

/** Print out all bank accounts sorted by name on STR. */
void printByName(PrintStream str) {
    for (Account account : names.values())
        account.print(str);
}
```

A Design Question: What would be an appropriate representation for keeping a record of all transactions (deposits and withdrawals) against each account?

Partial Implementations

- Besides interfaces (like List) and concrete types (like LinkedList), Java library provides abstract classes such as AbstractList.
- Idea is to take advantage of the fact that operations are related to each other.
- Example: once you know how to do get(k) and size() for an implementation of List, you can implement all the other methods needed for a *read-only* list (and its iterators).
- Now throw in add(k, x) and you have all you need for the additional operations of a growable list.
- Add set(k, x) and remove(k) and you can implement everything else.

Example: The `java.util.AbstractList` helper class

```
public abstract class AbstractList<Item> implements List<Item>
{
    /** Inherited from List */
    // public abstract int size();
    // public abstract Item get(int k);
    public boolean contains(Object x) {
        for (int i = 0; i < size(); i += 1) {
            if ((x == null && get(i) == null) ||
                (x != null && x.equals(get(i))))
                return true;
        }
        return false;
    }
    /* OPTIONAL: Throws exception; override to do more. */
    void add(int k, Item x) {
        throw new UnsupportedOperationException();
    }
}
```

Likewise for remove, set

Example, continued: AListIterator

```
// Continuing abstract class AbstractList<Item>:  
public Iterator<Item> iterator() { return listIterator(); }  
public ListIterator<Item> listIterator() {  
    return new AListIterator(this);  
}  
  
private static class AListIterator implements ListIterator<Item> {  
    AbstractList<Item> myList;  
    AListIterator(AbstractList<Item> L) { myList = L; }  
    /** Current position in our list. */  
    int where = 0;  
  
    public boolean hasNext() { return where < myList.size(); }  
    public Item next() { where += 1; return myList.get(where-1); }  
    public void add(Item x) { myList.add(where, x); where += 1; }  
    ... previous, remove, set, etc.  
}
```

...

Aside: Another way to do AListIterator

It's also possible to make the nested class non-static:

```
public Iterator<Item> iterator() { return listIterator(); }
public ListIterator<Item> listIterator() { return this.new AListIterator(); }

private class AListIterator implements ListIterator<Item> {
    /** Current position in our list. */
    int where = 0;

    public boolean hasNext() { return where < AbstractList.this.size(); }
    public Item next() { where += 1; return AbstractList.this.get(where-1); }
    public void add(Item x) { AbstractList.this.add(where, x); where += 1; }
    ... previous, remove, set, etc.
}
```

...

- Here, `AbstractList.this` means “the `AbstractList` I am attached to” and `X.new AListIterator` means “create a new `AListIterator` that is attached to `X`.”
- In this case you can abbreviate `this.new` as `new` and can leave off some `AbstractList.this` parts, since meaning is unambiguous.

Example: Using AbstractList

Problem: Want to create a *reversed view* of an existing List (same elements in reverse order). Operations on the original list affect the view, and vice-versa.

```
public class ReverseList<Item> extends AbstractList<Item> {  
    private final List<Item> L;  
  
    public ReverseList(List<Item> L) { this.L = L; }  
  
    public int size() { return L.size(); }  
  
    public Item get(int k) { return L.get(L.size()-k-1); }  
  
    public void add(int k, Item x) { L.add(L.size()-k, x); }  
  
    public Item set(int k, Item x) { return L.set(L.size()-k-1, x); }  
    public Item remove(int k) { return L.remove(L.size() - k - 1); }  
}
```

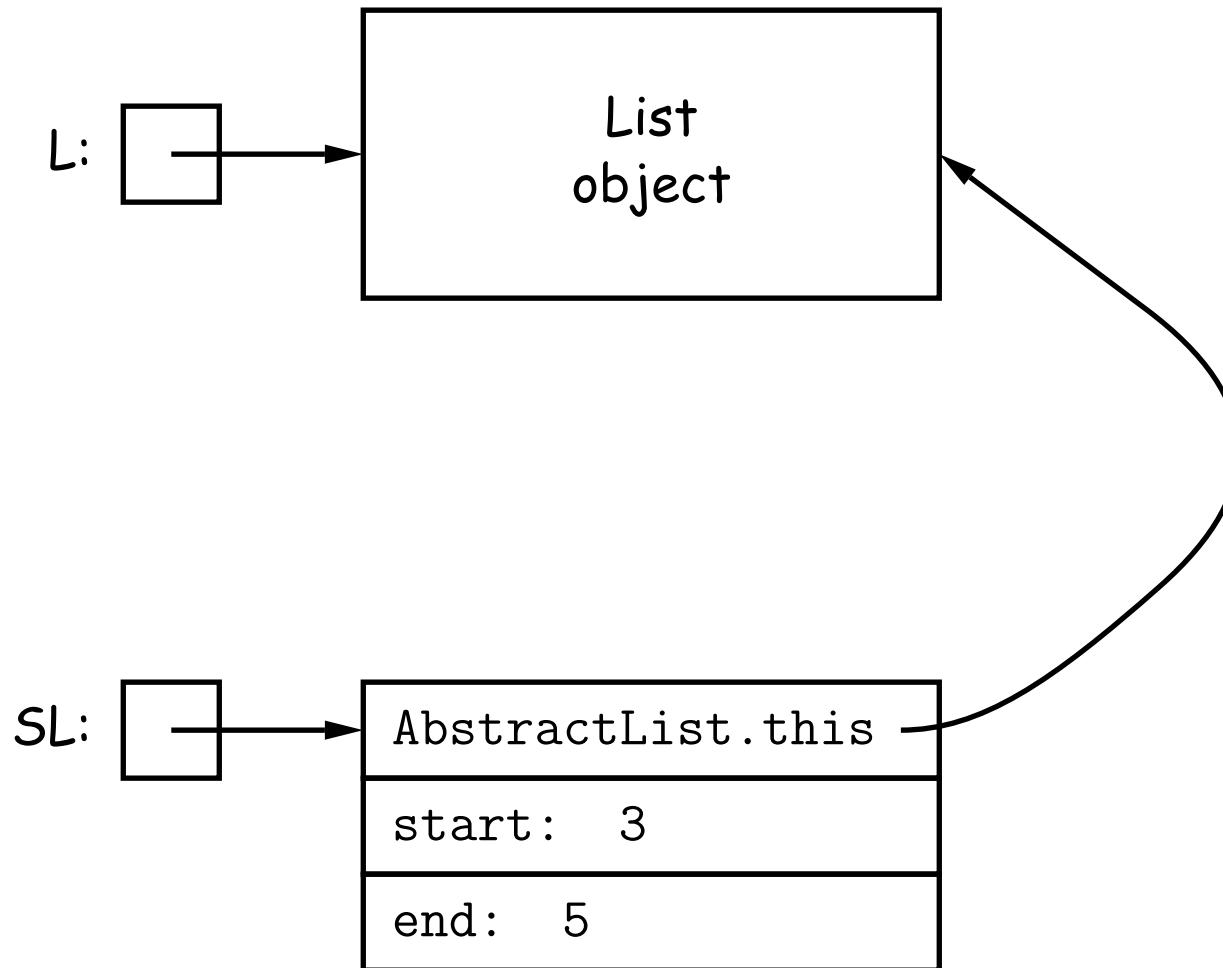
Getting a View: Sublists

Problem: `L.sublist(start, end)` is a List that gives a view of part of an existing list. Changes in one must affect the other. How?

```
// Continuation of class AbstractList. Error checks not shown.  
List<Item> sublist(int start, int end) {  
    return this.new Sublist(start, end);  
}  
  
private class Sublist extends AbstractList<Item> {  
    private int start, end;  
    Sublist(int start, int end) { obvious }  
  
    public int size() { return end-start; }  
    public Item get(int k) { return AbstractList.this.get(start+k); }  
  
    public void add(int k, Item x)  
    { AbstractList.this.add(start+k, x); end += 1; }  
    ...  
}
```

What Does a Sublist Look Like?

- Consider `SL = L.sublist(3, 5);`



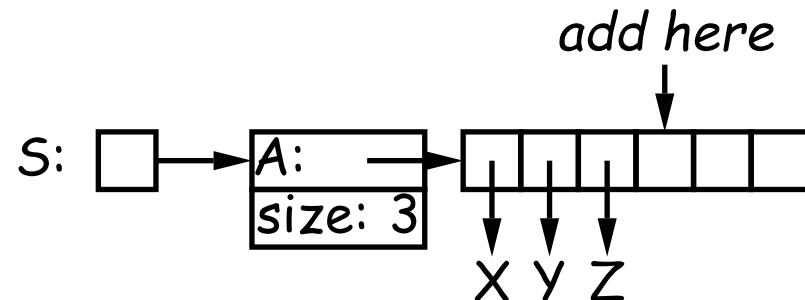
Arrays and Links

- Two main ways to represent a sequence: array and linked list
- In Java Library: ArrayList and Vector vs. LinkedList.
- Array:
 - Advantages: compact, fast ($\Theta(1)$) *random access* (indexing).
 - Disadvantages: insertion, deletion can be slow ($\Theta(N)$)
- Linked list:
 - Advantages: insertion, deletion fast once position found.
 - Disadvantages: space (link overhead), random access slow.

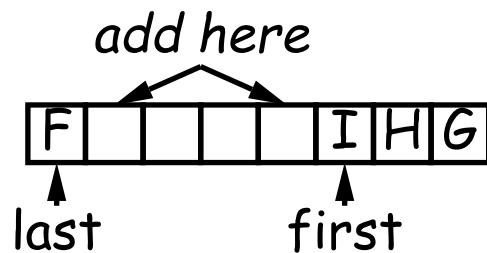
Implementing with Arrays

- Biggest problem using arrays is insertion/deletion in the *middle* of a list (must shove things over).
- Adding/deleting from ends can be made fast:
 - Double array size to grow; amortized cost constant (Lecture #15).
 - Growth at one end really easy; classical stack implementation:

```
S.push("X");  
S.push("Y");  
S.push("Z");
```



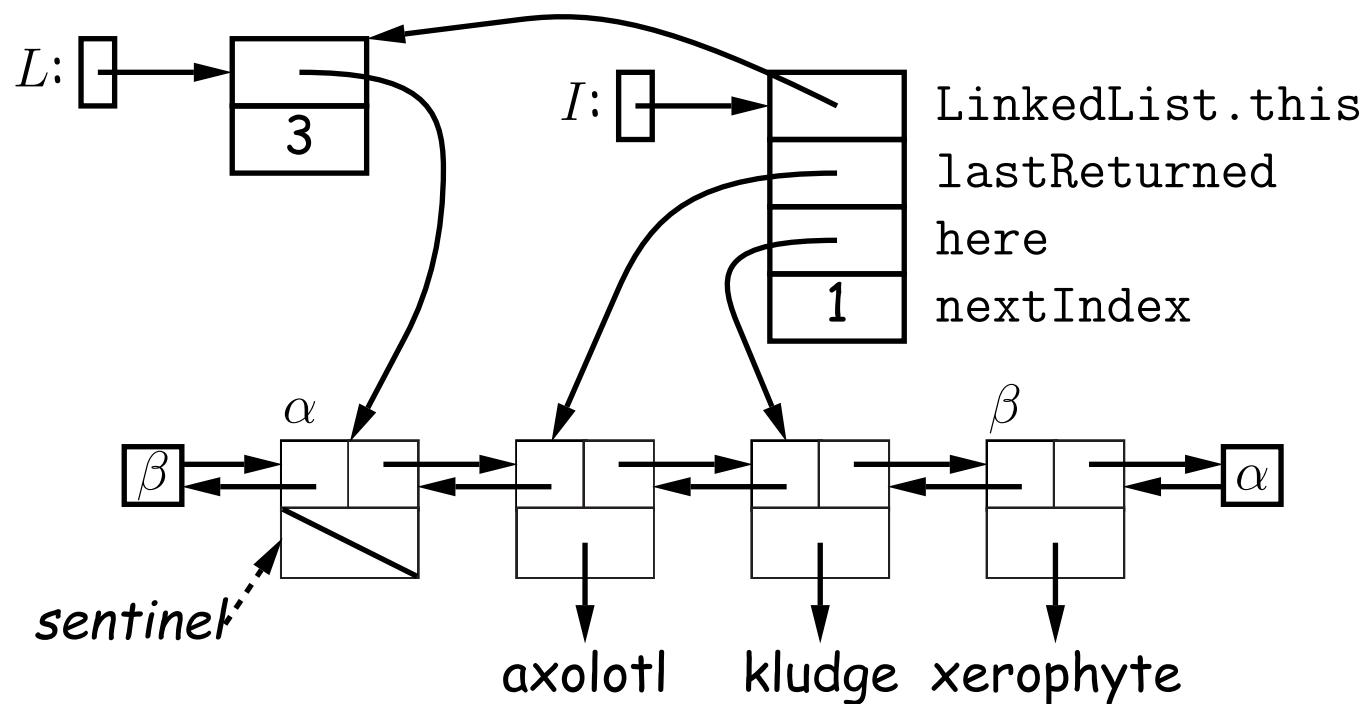
- To allow growth at either end, use *circular buffering*:



- Random access still fast.

Linking

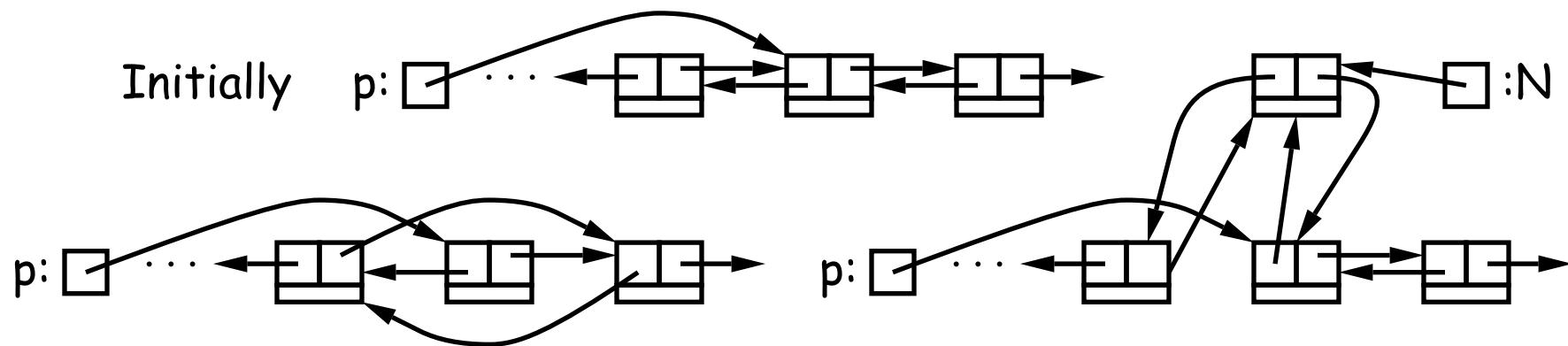
- Essentials of linking should now be familiar
- Used in Java LinkedList. One possible representation for linked list and an iterator object over it:



```
L = new LinkedList<String>();           I = L.listIterator();
L.add("axolotl");
L.add("kludge");
L.add("xerophyte");
```

Clever trick: Sentinels

- A *sentinel* is a dummy object containing no useful data except links.
- Used to eliminate special cases and to provide a fixed object to point to in order to access a data structure.
- Avoids special cases ('if' statements) by ensuring that the first and last item of a list always have (non-null) nodes—possibly sentinels—before and after them:
- `// To delete list node at p: // To add new node N before p:`
`p.next.prev = p.prev;` `N.prev = p.prev; N.next = p;`
`p.prev.next = p.next;` `p.prev.next = N;`
 `p.prev = N;`



Specialization

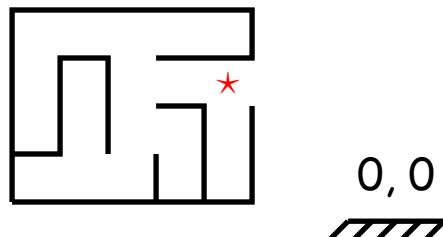
- Traditional special cases of general list:
 - **Stack**: Add and delete from one end (LIFO).
 - **Queue**: Add at end, delete from front (FIFO).
 - **Dequeue**: Add or delete at either end.
- All of these easily representable by either array (with circular buffering for queue or deque) or linked list.
- Java has the **List types**, which can act like any of these (although with non-traditional names for some of the operations).
- Also has `java.util.Stack`, a subtype of `List`, which gives traditional names ("push", "pop") to its operations. There is, however, no "stack" interface.

Stacks and Recursion

- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):  
    if isExit(start)  
        FOUND  
    else if (!isCrumb(start))  
        leave crumb at start;  
        for each square, x,  
            adjacent to start:  
                if legal(start,x) && !isCrumb(x)  
                    findExit(x)
```

Call: findExit(0,0)
Exit: (4, 2)



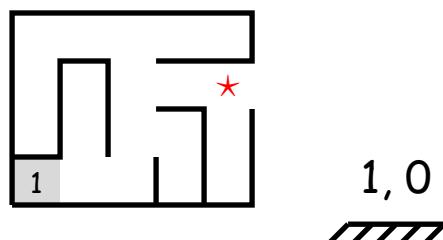
```
findExit(start):  
    S = new empty stack;  
    push start on S;  
    while S not empty:  
        pop S into start;  
        if isExit(start)  
            FOUND  
        else if (!isCrumb(start))  
            leave crumb at start;  
            for each square, x,  
                adjacent to start (in reverse):  
                    if legal(start,x) && !isCrumb(x)  
                        push x on S
```

Stacks and Recursion

- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):
    if isExit(start)
        FOUND
    else if (!isCrumb(start))
        leave crumb at start;
        for each square, x,
            adjacent to start:
                if legal(start,x) && !isCrumb(x)
                    findExit(x)
```

Call: findExit((0,0))
Exit: (4, 2)



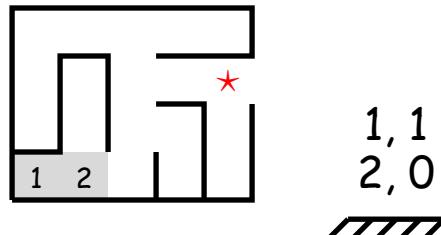
```
findExit(start):  
    S = new empty stack;  
    push start on S;  
    while S not empty:  
        pop S into start;  
        if isExit(start)  
            FOUND  
        else if (!isCrumb(start))  
            leave crumb at start;  
            for each square, x,  
                adjacent to start (in reverse):  
                if legal(start,x) && !isCrumb(x)  
                    push x on S
```

Stacks and Recursion

- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):  
    if isExit(start)  
        FOUND  
    else if (!isCrumb(start))  
        leave crumb at start;  
        for each square, x,  
            adjacent to start:  
                if legal(start,x) && !isCrumb(x)  
                    findExit(x)
```

Call: findExit(0,0)
Exit: (4, 2)



```
findExit(start):  
    S = new empty stack;  
    push start on S;  
    while S not empty:  
        pop S into start;  
        if isExit(start)  
            FOUND  
        else if (!isCrumb(start))  
            leave crumb at start;  
            for each square, x,  
                adjacent to start (in reverse):  
                    if legal(start,x) && !isCrumb(x)  
                        push x on S
```

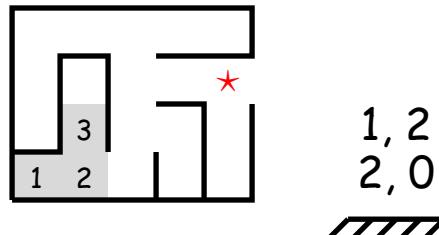
Stacks and Recursion

- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):  
    if isExit(start)  
        FOUND  
    else if (!isCrumb(start))  
        leave crumb at start;  
        for each square, x,  
            adjacent to start:  
                if legal(start,x) && !isCrumb(x)  
                    findExit(x)
```

Call: `findExit(0,0)`

Exit: (4, 2)



```
findExit(start):  
    S = new empty stack;  
    push start on S;  
    while S not empty:  
        pop S into start;  
        if isExit(start)  
            FOUND  
        else if (!isCrumb(start))  
            leave crumb at start;  
            for each square, x,  
                adjacent to start (in reverse):  
                    if legal(start,x) && !isCrumb(x)  
                        push x on S
```

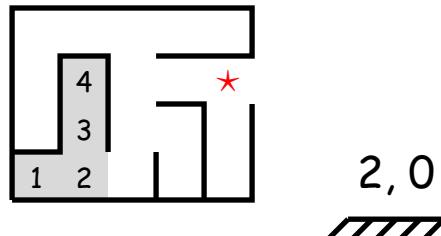
Stacks and Recursion

- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):  
    if isExit(start)  
        FOUND  
    else if (!isCrumb(start))  
        leave crumb at start;  
        for each square, x,  
            adjacent to start:  
                if legal(start,x) && !isCrumb(x)  
                    findExit(x)
```

Call: `findExit(0,0)`

Exit: (4, 2)



```
findExit(start):  
    S = new empty stack;  
    push start on S;  
    while S not empty:  
        pop S into start;  
        if isExit(start)  
            FOUND  
        else if (!isCrumb(start))  
            leave crumb at start;  
            for each square, x,  
                adjacent to start (in reverse):  
                    if legal(start,x) && !isCrumb(x)  
                        push x on S
```

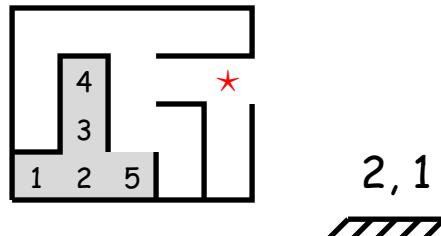
Stacks and Recursion

- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):  
    if isExit(start)  
        FOUND  
    else if (!isCrumb(start))  
        leave crumb at start;  
        for each square, x,  
            adjacent to start:  
                if legal(start,x) && !isCrumb(x)  
                    findExit(x)
```

Call: `findExit(0,0)`

Exit: (4, 2)



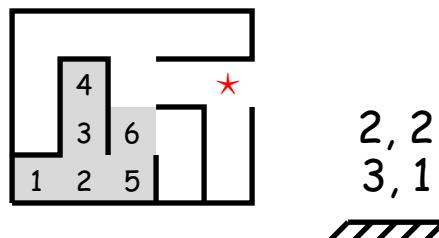
```
findExit(start):  
    S = new empty stack;  
    push start on S;  
    while S not empty:  
        pop S into start;  
        if isExit(start)  
            FOUND  
        else if (!isCrumb(start))  
            leave crumb at start;  
            for each square, x,  
                adjacent to start (in reverse):  
                    if legal(start,x) && !isCrumb(x)  
                        push x on S
```

Stacks and Recursion

- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):  
    if isExit(start)  
        FOUND  
    else if (!isCrumb(start))  
        leave crumb at start;  
        for each square, x,  
            adjacent to start:  
                if legal(start,x) && !isCrumb(x)  
                    findExit(x)
```

Call: `findExit(0,0)`
Exit: (4, 2)



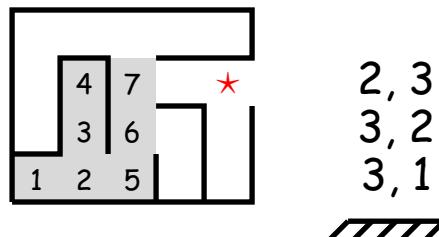
```
findExit(start):  
    S = new empty stack;  
    push start on S;  
    while S not empty:  
        pop S into start;  
        if isExit(start)  
            FOUND  
        else if (!isCrumb(start))  
            leave crumb at start;  
            for each square, x,  
                adjacent to start (in reverse):  
                    if legal(start,x) && !isCrumb(x)  
                        push x on S
```

Stacks and Recursion

- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):  
    if isExit(start)  
        FOUND  
    else if (!isCrumb(start))  
        leave crumb at start;  
        for each square, x,  
            adjacent to start:  
                if legal(start,x) && !isCrumb(x)  
                    findExit(x)
```

Call: `findExit(0,0)`
Exit: (4, 2)



```
findExit(start):  
    S = new empty stack;  
    push start on S;  
    while S not empty:  
        pop S into start;  
        if isExit(start)  
            FOUND  
        else if (!isCrumb(start))  
            leave crumb at start;  
            for each square, x,  
                adjacent to start (in reverse):  
                    if legal(start,x) && !isCrumb(x)  
                        push x on S
```

Stacks and Recursion

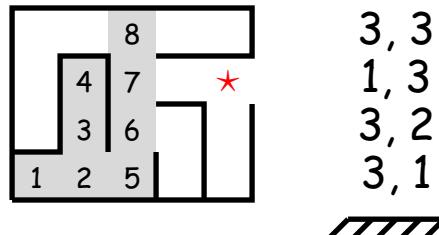
- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):
```

```
    if isExit(start)
        FOUND
    else if (!isCrumb(start))
        leave crumb at start;
        for each square, x,
            adjacent to start:
                if legal(start,x) && !isCrumb(x)
                    findExit(x)
```

```
Call: findExit((0,0))
```

```
Exit: (4, 2)
```



```
findExit(start):
```

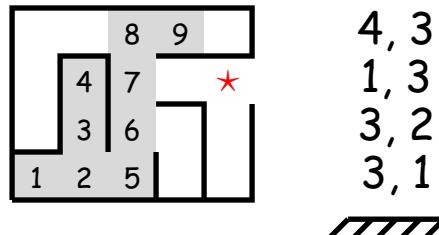
```
    S = new empty stack;
    push start on S;
    while S not empty:
        pop S into start;
        if isExit(start)
            FOUND
        else if (!isCrumb(start))
            leave crumb at start;
            for each square, x,
                adjacent to start (in reverse):
                    if legal(start,x) && !isCrumb(x)
                        push x on S
```

Stacks and Recursion

- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):  
    if isExit(start)  
        FOUND  
    else if (!isCrumb(start))  
        leave crumb at start;  
        for each square, x,  
            adjacent to start:  
                if legal(start,x) && !isCrumb(x)  
                    findExit(x)
```

Call: findExit((0,0))
Exit: (4, 2)



```
findExit(start):  
    S = new empty stack;  
    push start on S;  
    while S not empty:  
        pop S into start;  
        if isExit(start)  
            FOUND  
        else if (!isCrumb(start))  
            leave crumb at start;  
            for each square, x,  
                adjacent to start (in reverse):  
                    if legal(start,x) && !isCrumb(x)  
                        push x on S
```

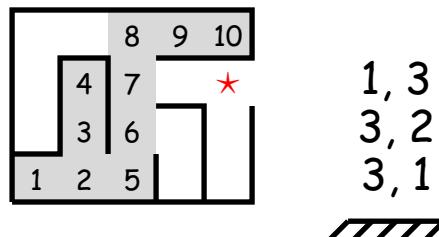
Stacks and Recursion

- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):  
    if isExit(start)  
        FOUND  
    else if (!isCrumb(start))  
        leave crumb at start;  
        for each square, x,  
            adjacent to start:  
                if legal(start,x) && !isCrumb(x)  
                    findExit(x)
```

Call: `findExit(0,0)`

Exit: (4, 2)



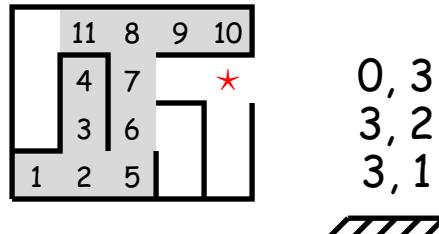
```
findExit(start):  
    S = new empty stack;  
    push start on S;  
    while S not empty:  
        pop S into start;  
        if isExit(start)  
            FOUND  
        else if (!isCrumb(start))  
            leave crumb at start;  
            for each square, x,  
                adjacent to start (in reverse):  
                    if legal(start,x) && !isCrumb(x)  
                        push x on S
```

Stacks and Recursion

- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):  
    if isExit(start)  
        FOUND  
    else if (!isCrumb(start))  
        leave crumb at start;  
        for each square, x,  
            adjacent to start:  
                if legal(start,x) && !isCrumb(x)  
                    findExit(x)
```

Call: `findExit(0,0)`
Exit: (4, 2)



```
findExit(start):  
    S = new empty stack;  
    push start on S;  
    while S not empty:  
        pop S into start;  
        if isExit(start)  
            FOUND  
        else if (!isCrumb(start))  
            leave crumb at start;  
            for each square, x,  
                adjacent to start (in reverse):  
                    if legal(start,x) && !isCrumb(x)  
                        push x on S
```

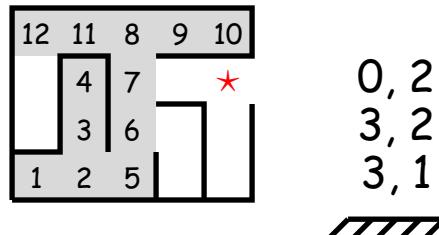
Stacks and Recursion

- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):  
    if isExit(start)  
        FOUND  
    else if (!isCrumb(start))  
        leave crumb at start;  
        for each square, x,  
            adjacent to start:  
                if legal(start,x) && !isCrumb(x)  
                    findExit(x)
```

Call: findExit(0,0)

Exit: (4, 2)



```
findExit(start):  
    S = new empty stack;  
    push start on S;  
    while S not empty:  
        pop S into start;  
        if isExit(start)  
            FOUND  
        else if (!isCrumb(start))  
            leave crumb at start;  
            for each square, x,  
                adjacent to start (in reverse):  
                    if legal(start,x) && !isCrumb(x)  
                        push x on S
```

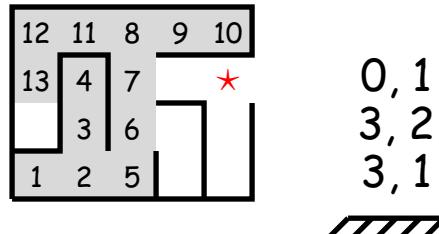
Stacks and Recursion

- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):  
    if isExit(start)  
        FOUND  
    else if (!isCrumb(start))  
        leave crumb at start;  
        for each square, x,  
            adjacent to start:  
                if legal(start,x) && !isCrumb(x)  
                    findExit(x)
```

Call: findExit(0,0)

Exit: (4, 2)



```
findExit(start):  
    S = new empty stack;  
    push start on S;  
    while S not empty:  
        pop S into start;  
        if isExit(start)  
            FOUND  
        else if (!isCrumb(start))  
            leave crumb at start;  
            for each square, x,  
                adjacent to start (in reverse):  
                    if legal(start,x) && !isCrumb(x)  
                        push x on S
```

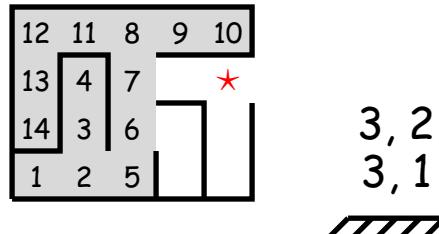
Stacks and Recursion

- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):  
    if isExit(start)  
        FOUND  
    else if (!isCrumb(start))  
        leave crumb at start;  
        for each square, x,  
            adjacent to start:  
                if legal(start,x) && !isCrumb(x)  
                    findExit(x)
```

Call: `findExit(0,0)`

Exit: (4, 2)



```
findExit(start):  
    S = new empty stack;  
    push start on S;  
    while S not empty:  
        pop S into start;  
        if isExit(start)  
            FOUND  
        else if (!isCrumb(start))  
            leave crumb at start;  
            for each square, x,  
                adjacent to start (in reverse):  
                    if legal(start,x) && !isCrumb(x)  
                        push x on S
```

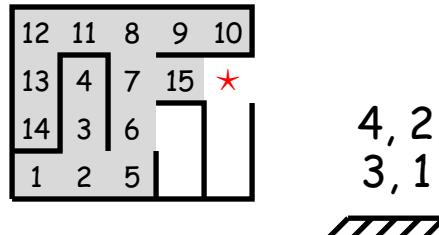
Stacks and Recursion

- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):  
    if isExit(start)  
        FOUND  
    else if (!isCrumb(start))  
        leave crumb at start;  
        for each square, x,  
            adjacent to start:  
                if legal(start,x) && !isCrumb(x)  
                    findExit(x)
```

Call: `findExit(0,0)`

Exit: (4, 2)



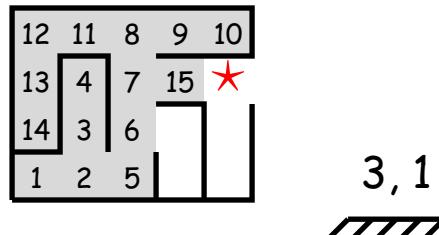
```
findExit(start):  
    S = new empty stack;  
    push start on S;  
    while S not empty:  
        pop S into start;  
        if isExit(start)  
            FOUND  
        else if (!isCrumb(start))  
            leave crumb at start;  
            for each square, x,  
                adjacent to start (in reverse):  
                    if legal(start,x) && !isCrumb(x)  
                        push x on S
```

Stacks and Recursion

- Stacks are related to recursion. In fact, can convert any recursive algorithm to stack-based (however, generally with no great performance benefit):
 - Calls become “push current variables and parameters, set parameters to new values, and loop.”
 - Return becomes “pop to restore variables and parameters.”

```
findExit(start):  
    if isExit(start)  
        FOUND  
    else if (!isCrumb(start))  
        leave crumb at start;  
        for each square, x,  
            adjacent to start:  
                if legal(start,x) && !isCrumb(x)  
                    findExit(x)
```

Call: `findExit(0,0)`
Exit: (4, 2)



```
findExit(start):  
    S = new empty stack;  
    push start on S;  
    while S not empty:  
        pop S into start;  
        if isExit(start)  
            FOUND  
        else if (!isCrumb(start))  
            leave crumb at start;  
            for each square, x,  
                adjacent to start (in reverse):  
                    if legal(start,x) && !isCrumb(x)  
                        push x on S
```

Design Choices: Extension, Delegation, Adaptation

- The standard `java.util.Stack` type *extends* `Vector`:

```
class Stack<Item> extends Vector<Item> { void push(Item x) { add(x); } ... }
```

- Could instead have *delegated* to a field:

```
class ArrayStack<Item> {  
    private ArrayList<Item> repl = new ArrayList<Item>();  
    void push(Item x) { repl.add(x); } ...  
}
```

- Or, could generalize, and define an *adapter*: a class used to make objects of one kind behave as another:

```
public class StackAdapter<Item> {  
    private List repl;  
    /** A stack that uses REPL for its storage. */  
    public StackAdapter(List<Item> repl) { this.repl = repl; }  
    public void push(Item x) { repl.add(x); } ...  
}
```

```
class ArrayStack<Item> extends StackAdapter<Item> {  
    ArrayStack() { super(new ArrayList<Item>()); }  
}
```

CS61B Lecture #20: Trees

A Recursive Structure

- Trees naturally represent recursively defined, hierarchical objects with more than one recursive subpart for each instance.
- Common examples: expressions, sentences.
 - Expressions have definitions such as "an expression consists of a literal or two expressions separated by an operator."
- Also describe search structures in which we recursively divide a set into multiple disjoint subsets.

Formal Definitions

- Trees come in a variety of flavors, all defined recursively:
 - **61A style:** A tree consists of a *label* value and zero or more *branches* (or *children*), each of them a tree.
 - **61A style, alternative definition:** A tree is a set of *nodes* (or *vertices*), each of which has a label value and one or more *child nodes*, such that no node descends (directly or indirectly) from itself. A node is the *parent* of its children.
 - **Positional trees:** A tree is either *empty* or consists of a node containing a label value and an indexed sequence of zero or more children, each a positional tree. If every node has two positions, we have a *binary tree* and the children are its *left and right subtrees*. Again, nodes are the parents of their non-empty children.
 - We'll see other varieties when considering graphs.

Tree Characteristics (I)

- The *root* of a tree is a non-empty node with no parent in that tree (its parent might be in some larger tree that contains that tree as a subtree). Thus, every node is the root of a (sub)tree.
- The *order*, *arity*, or *degree* of a node (tree) is its number (maximum number) of children.
- The nodes of a *k-ary tree* each have at most k children.
- A *leaf* node has no children (no non-empty children in the case of positional trees).

Tree Characteristics (II)

- The *height* of a node in a tree is the largest distance to a leaf. That is, a leaf has height 0 and a non-empty tree's height is one more than the maximum height of its children. The height of a tree is the height of its root.
- The *depth* of a node in a tree is the distance to the root of that tree. That is, in a tree whose root is R , R itself has depth 0 in R , and if node $S \neq R$ is in the tree with root R , then its depth is one greater than its parent's.

A Tree Type, 61A Style

```
/** An immutable Tree whose labels are of type LABEL. */
public class Tree<Label> {

    @SuppressWarnings("unchecked") // Explained in a later lecture
    public Tree(Label label, Tree<Label>... children) {
        _label = label;
        _kids = new ArrayList<>(Arrays.asList(children));
    }

    public int arity() { return _kids.size(); }

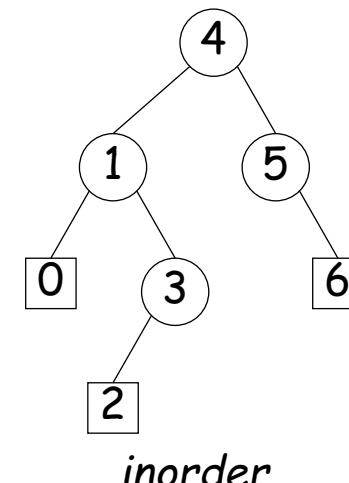
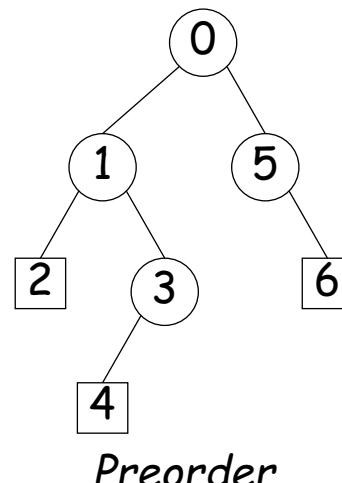
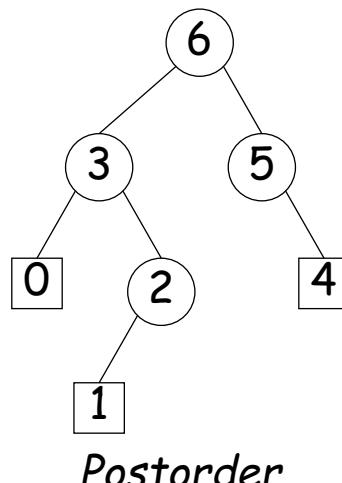
    public Label label() { return _label; }

    public Tree<Label> child(int k) { return _kids.get(k); }

    private Label _label;
    private ArrayList<Tree<Label>> _kids;
}
```

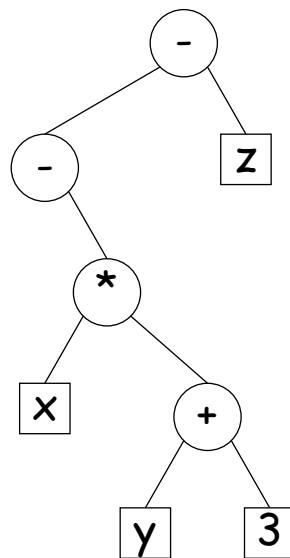
Fundamental Operation: Traversal

- *Traversing a tree* means enumerating (some subset of) its nodes.
- Typically done recursively, because that is natural description.
- As nodes are enumerated, we say they are *visited*.
- Three basic orders for enumeration (+ variations):
 - **Preorder**: visit node, traverse its children.
 - **Postorder**: traverse children, visit node.
 - **Inorder**: traverse first child, visit node, traverse second child (binary trees only).



Preorder Traversal and Prefix Expressions

Problem: Convert



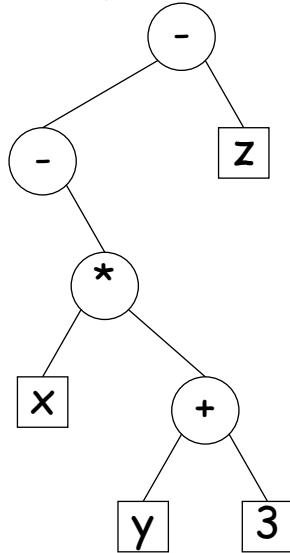
into

(- (- (* x (+ y 3))) z)

```
static String toLisp(Tree<String> T) {  
    if (T.arity() == 0) return T.label();  
    else {  
        String R; R = "(" + T.label();  
        for (int i = 0; i < T.arity(); i += 1)  
            R += " " + toLisp(T.child(i));  
        return R + ")";  
    }  
}
```

Inorder Traversal and Infix Expressions

Problem: Convert



into

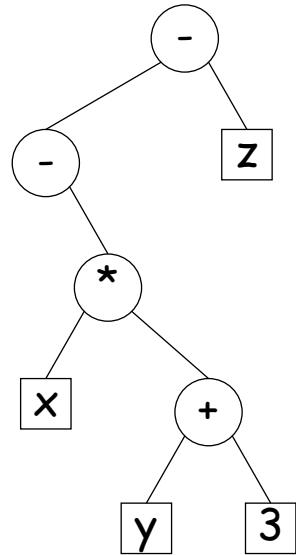
$((-(x*(y+3)))-z)$

To think about: how to get rid of all those parentheses.

```
static String toInfix(Tree<String> T) {  
    if (T.arity() == 0) {  
        return T.label();  
    } else if (T.arity() == 1) {  
        return "(" + T.label() + toInfix(T.child(0)) + ")";  
    } else {  
        return "(" + toInfix(T.child(0)) + T.label() + toInfix(T.child(1)) + ")";  
    }  
}
```

Postorder Traversal and Postfix Expressions

Problem: Convert



into

x y 3 +:2 *:2 -:1 z -:2

```
static String toPolish(Tree<String> T) {
    String R; R = "";
    for (int i = 0; i < T.arity(); i += 1)
        R += toPolish(T.child(i)) + " ";
    return R + String.format("%s:%d", T.label(), T.arity());
}
```

A General Traversal: The Visitor Pattern

```
void preorderTraverse(Tree<Label> T, Consumer<Tree<Label>> visit)
{
    if (T != null) {
        visit.accept(T);
        for (int i = 0; i < T.arity(); i += 1)
            preorderTraverse(T.child(i), visit);
    }
}
```

- `java.util.function.Consumer<AType>` is a library interface that works as a function-like type with one void method, `accept`, which takes an argument of type `AType`.
- Now, using Java 8 lambda syntax, I can print all labels in the tree in preorder with:

```
preorderTraverse(myTree, T -> System.out.print(T.label() + " "));
```

Iterative Depth-First Traversals

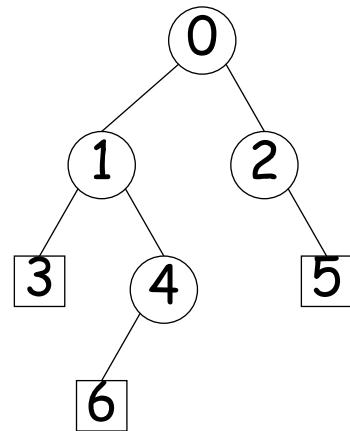
- Tree recursion conceals data: a *stack* of nodes (all the T arguments) and a little extra information. Can make the data explicit:

```
void preorderTraverse2(Tree<Label> T, Consumer<Tree<Label>> visit) {  
    Stack<Tree<Label>> work = new Stack<>();  
    work.push(T);  
    while (!work.isEmpty()) {  
        Tree<Label> node = work.pop();  
        visit.accept(node);  
        for (int i = node.arity()-1; i >= 0; i -= 1)  
            work.push(node.child(i)); // Why backward?  
    }  
}
```

- This traversal takes the same $\Theta(\cdot)$ time as doing it recursively, and also the same $\Theta(\cdot)$ space.
- That is, we have substituted an explicit stack data structure (work) for Java's built-in execution stack (which handles function calls).

Level-Order (Breadth-First) Traversal

Problem: Traverse all nodes at depth 0, then depth 1, etc:



Breadth-First Traversal Implemented

A simple modification to iterative depth-first traversal gives breadth-first traversal. Just change the (LIFO) stack to a (FIFO) queue:

```
void breadthFirstTraverse(Tree<Label> T, Consumer<Tree<Label>> visit) {  
    ArrayDeque<Tree<Label>> work = new ArrayDeque<>(); // (Changed)  
    work.push(T);  
    while (!work.isEmpty()) {  
        Tree<Label> node = work.remove(); // (Changed)  
        if (node != null) {  
            visit.accept(node);  
            for (int i = 0; i < node.arity(); i += 1) // (Changed)  
                work.push(node.child(i));  
        }  
    }  
}
```

Times

- The traversal algorithms have roughly the form of the boom example in §1.3.3 of *Data Structures*—an exponential algorithm.
- However, the role of M in that algorithm is played by the *height* of the tree, not the number of nodes.
- In fact, easy to see that tree traversal is *linear*: $\Theta(N)$, where N is the # of nodes: Form of the algorithm implies that there is one visit at the root, and then one visit for every *edge* in the tree. Since every node but the root has exactly one parent, and the root has none, *must be $N - 1$ edges in any non-empty tree*.
- In positional tree, is also one recursive call for each empty tree, but # of empty trees can be no greater than kN , where k is arity.
- For k -ary tree (max # children is k), $h + 1 \leq N \leq \frac{k^{h+1} - 1}{k - 1}$, where h is height.
- So $h \in \Omega(\log_k N) = \Omega(\lg N)$ and $h \in O(N)$.
- Many tree algorithms look at one child only. For them, worst-case time is proportional to the *height* of the tree— $\Theta(\lg N)$ —assuming that tree is *bushy*—each level has about as many nodes as possible.

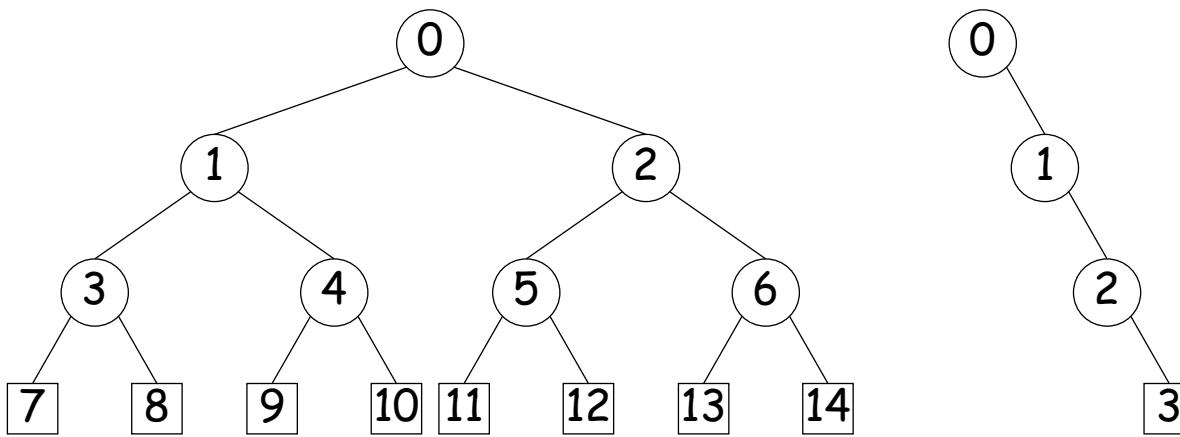
Recursive Breadth-First Traversal: Iterative Deepening

- Previous breadth-first traversal used space proportional to the *width* of the tree, which is $\Theta(N)$ for bushy trees, whereas depth-first traversal takes $\lg N$ space on bushy trees.
- Can we get breadth-first traversal in $\lg N$ space and $\Theta(N)$ time on bushy trees?
- For each level, k , of the tree from 0 to *lev*, call `doLevel(T, k)`:

```
void doLevel(Tree T, int lev) {  
    if (lev == 0)  
        visit T  
    else  
        for each non-null child, C, of T {  
            doLevel(C, lev-1);  
        }  
}
```

- So we do breadth-first traversal by repeated (truncated) depth-first traversals: *iterative deepening*.
- In `doLevel(T, k)`, we skip (i.e., traverse but don't visit) the nodes before level k , and then visit at level k , but not their children.

Iterative Deepening Time?



- Let h be height, N be # of nodes.
- Count # edges traversed (i.e, # of calls, not counting null nodes).
- First (full) tree: 1 for level 0, 3 for level 1, 7 for level 2, 15 for level 3.
- Or in general $(2^1 - 1) + (2^2 - 1) + \dots + (2^{h+1} - 1) = 2^{h+2} - 2 - (h+1) \in \Theta(N)$, since $N = 2^{h+1} - 1$ for this tree.
- Second (right leaning) tree: 1 for level 0, 2 for level 2, 3 for level 3.
- Or in general $(h+1)(h+2)/2 = N(N+1)/2 \in \Theta(N^2)$, since $N = h+1$ for this kind of tree.

Iterators for Trees

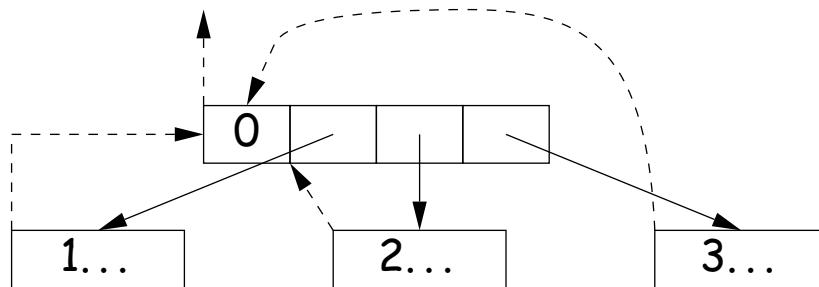
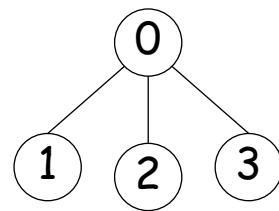
- Frankly, iterators are not terribly convenient on trees.
- But can use ideas from iterative methods.

```
class PreorderTreeIterator<Label> implements Iterator<Label> {  
    private Stack<Tree<Label>> s = new Stack<Tree<Label>>();  
  
    public PreorderTreeIterator(Tree<Label> T) { s.push(T); }  
  
    public boolean hasNext() { return !s.isEmpty(); }  
    public T next() {  
        Tree<Label> result = s.pop();  
        for (int i = result.arity()-1; i >= 0; i -= 1)  
            s.push(result.child(i));  
        return result.label();  
    }  
}
```

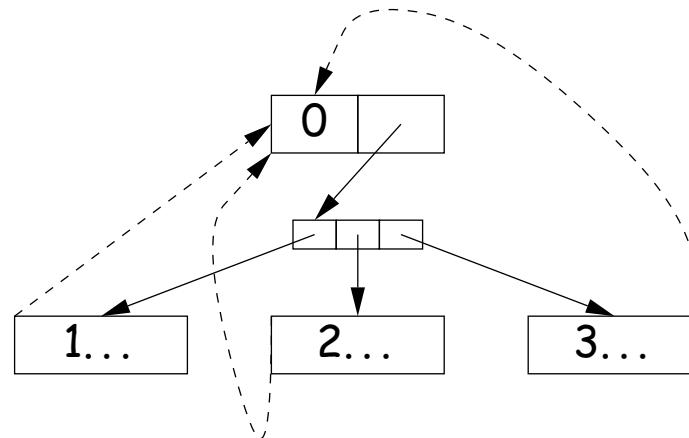
Example: (what do I have to add to class Tree first?)

```
for (String label : aTree) System.out.print(label + " ");
```

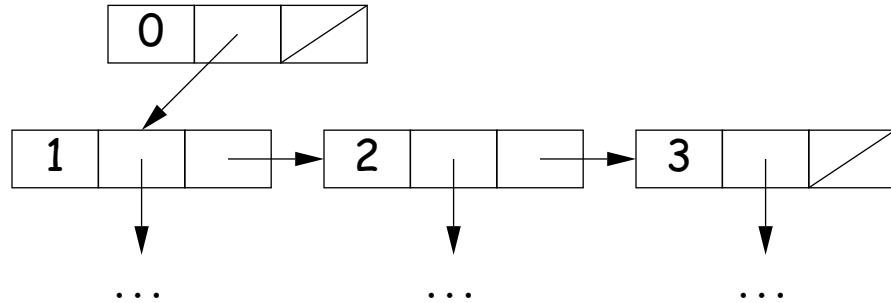
Tree Representation



(a) Embedded child pointers
(+ optional parent pointers)



(b) Array of child pointers
(+ optional parent pointers)



(c) child/sibling pointers

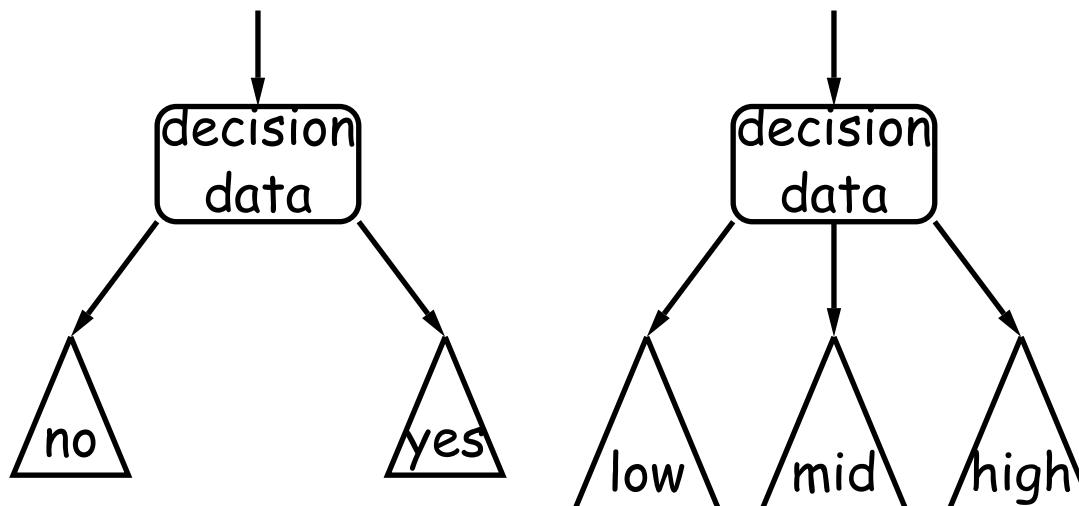


(d) breadth-first array
(complete trees)

CS61B Lecture #21: Tree Searching

Divide and Conquer

- Much (most?) computation is devoted to finding things in response to various forms of query.
- Linear search for response can be expensive, especially when data set is too large for primary memory.
- Preferable to have criteria for *dividing* data to be searched into pieces recursively
- We saw the figure for $\lg N$ algorithms: at 1 μ sec per comparison, could process 10^{300000} items in 1 sec.
- Tree is a natural framework for the representation:



Binary Search Trees

Binary Search Property:

- Tree nodes contain *keys*, and possibly other data.
- All nodes in left subtree of node have *smaller* keys.
- All nodes in right subtree of node have *larger* keys.
- “Smaller” means any complete transitive, anti-symmetric ordering on keys:
 - exactly one of $x \prec y$ and $y \prec x$ true.
 - $x \prec y$ and $y \prec z$ imply $x \prec z$.
 - (To simplify, won’t allow duplicate keys this semester).
- E.g., in dictionary database, node label would be (word, definition): word is the key.
- For concreteness here, we’ll just use the standard Java convention of calling .compareTo.

A Binary Search Type

- Here, we'll use the following simple binary search tree type. Ignore all the style violations, please.

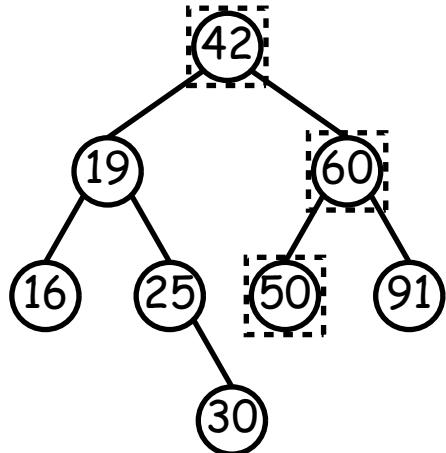
```
/** A node of a binary search tree associating a value of type VALUE
 * with a key of type KEY. (Thus, the labels in this tree are
 * key/value pairs.) */
class BST<Key extends Comparable<Key>, Value> {
    Key key;
    Value value;
    BST<Key, Value> left, right;

    BST(Key key0, Value value0,
        BST<Key, Value> left0, BST<Key, Value> right0) {
        Body left to the reader.
    }
    BST(Key key0, Value value0) {
        this(key0, value0, null, null);
    }
}
```

- (Ignore the Key extends Comparable<Key> stuff for now. It just says that keys (of type Key) can be compared to each other.)

Finding

- Searching for 50 and 49:

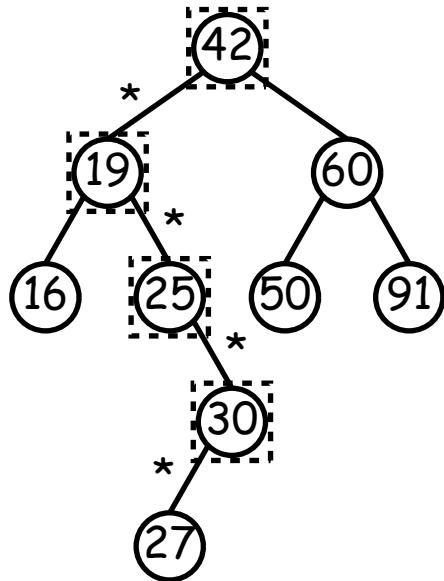


```
/** Return node in T containing L. Null if none. */
static <Key extends Comparable<Key>, Value>
BST<Key, Value> find(BST<Key, Value> T, Key L) {
    if (T == null)
        return T;
    if (L.compareTo(T.key) == 0)
        return T;
    else if (L.compareTo(T.key) < 0)
        return find(T.left, L);
    else
        return find(T.right, L);
}
```

- Dashed boxes show which node labels we look at.
- Number of nodes examined is proportional to height of tree.

Inserting

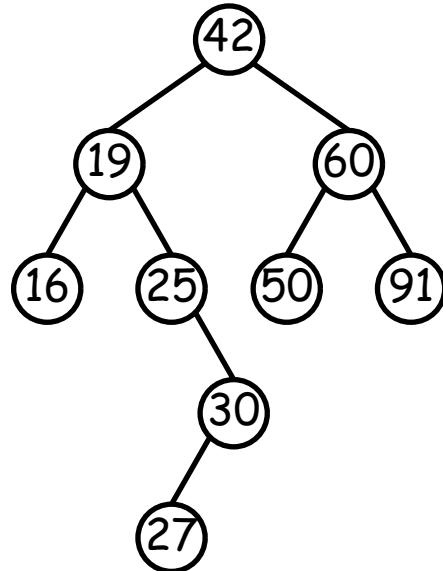
- Inserting 27



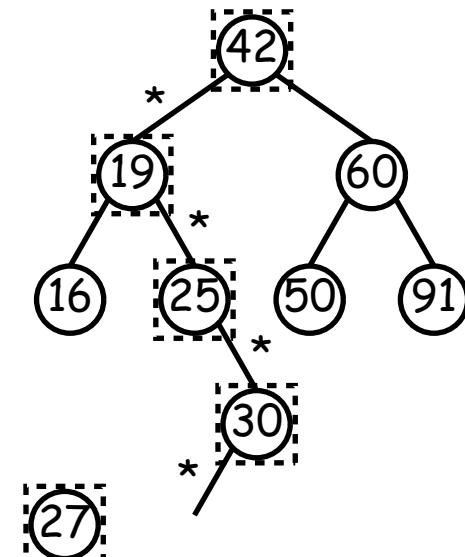
```
/** Insert V in T with key K, replacing existing
 * value if present. Return the modified tree.*/
static <Key extends Comparable<Key>, Value>
BST<Key, Value> insert(BST<Key, Value> T,
                      Key K, Value V) {
    if (T == null)
        return new BST(K, V);
    if (K.compareTo(T.key) == 0)
        T.value = V;
    else if (K.compareTo(T.key) < 0)
        T.left = insert(T.left, K, V);
    else
        T.right = insert(T.right, K, V);
    return T;
}
```

- Starred edges are set (to themselves, unless initially null).
- Again, time proportional to height.

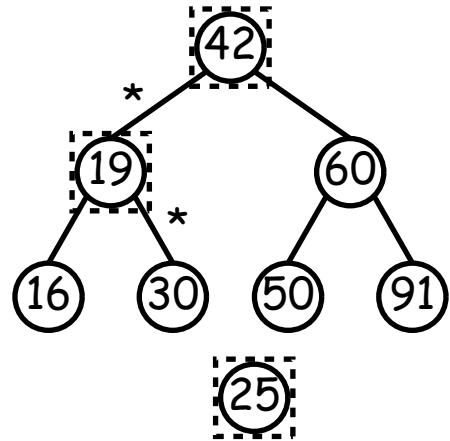
Deletion



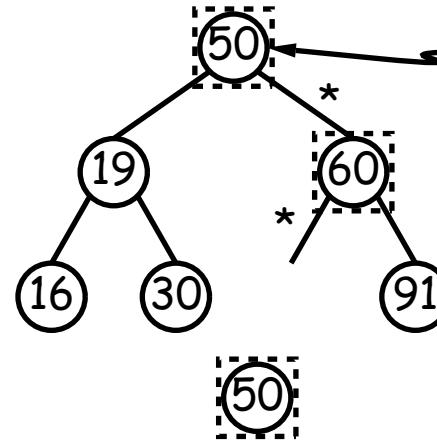
Initial



Remove 27



Remove 25



Remove 42

Deletion Algorithm

```
/** Remove K from T, and return the new tree. */
static <Key extends Comparable<Key>, Value>
BST<Key, Value> remove(BST T, Key K) {

    if (T == null)
        return null;
    if (K.compareTo(T.key) == 0) {
        if (T.left == null)
            return T.right;
        else if (T.right == null)
            return T.left;
        else {
            BST<Key, Value> smallest = minNode(T.right); // ??
            T.value = smallest.value;
            T.key = smallest.key;
            T.right = remove(T.right, smallest.key);
        }
    }
    else if (K.compareTo(T.key) < 0)
        T.left = remove(T.left, K);
    else
        T.right = remove(T.right, K);
    return T;
}
```

Minimum Node in a Subtree

```
/** Return the node containing the minimal key T.  
 * T must not be null. */  
public static <Key extends Comparable<Key>, Value>  
BST<Key, Value> minNode(BST<Key, Value> T) {  
  
    // ???  
  
    return T;  
}
```

- This is a *generic function*. As before, let's ignore all the "<Key extends Comparable...>" stuff for now.

Minimum Node in a Subtree

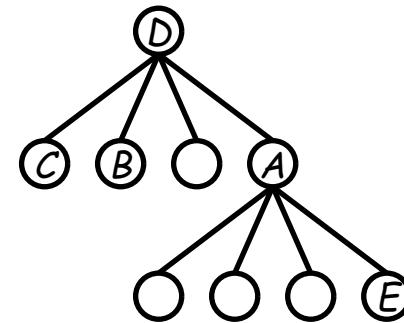
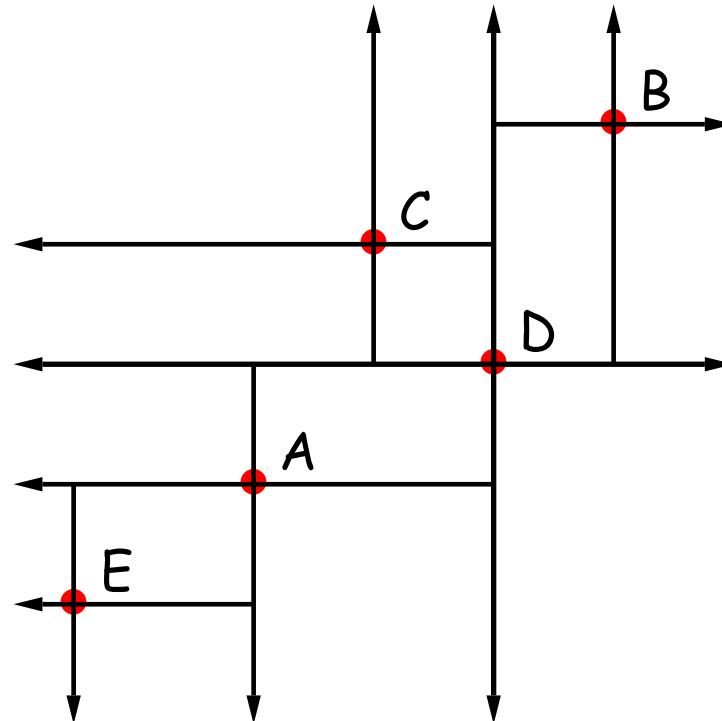
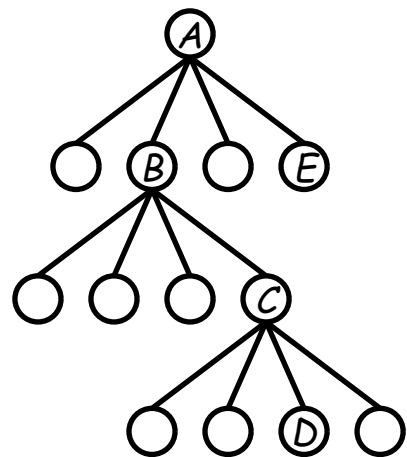
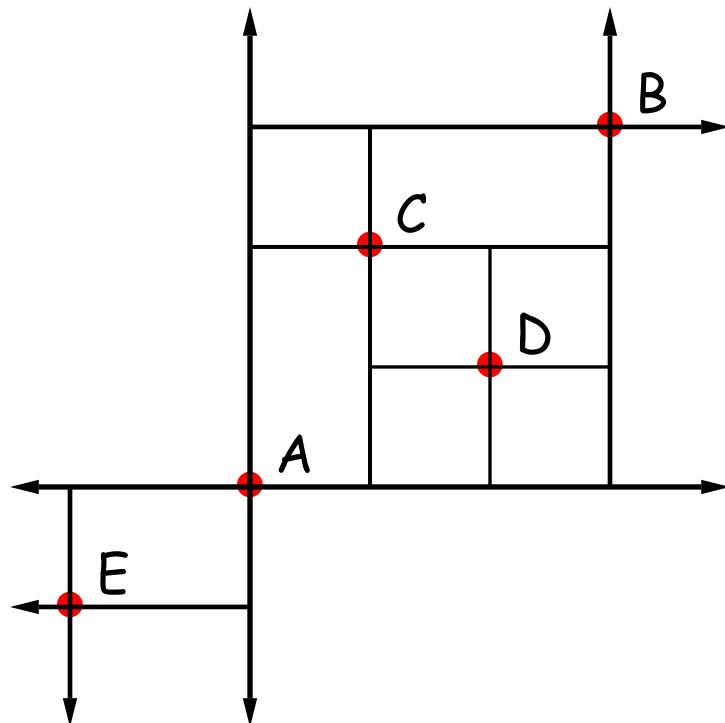
```
/** Return the node containing the minimal key T.  
 * T must not be null. */  
public static <Key extends Comparable<Key>, Value>  
BST<Key, Value> minNode(BST<Key, Value> T) {  
  
    while (T.left != null) {  
        T = T.left;  
    }  
  
    return T;  
}
```

- This is a *generic function*. As before, let's ignore all the "<Key extends Comparable...>" stuff for now.

More Than Two Choices: Quadtrees

- Want to **index** information about 2D locations so that items can be retrieved by position. But how to compare positions "binarily?"
- **Quadtrees** do so using the same standard data-structuring trick as BSTs—**Divide and Conquer**—but with more subtrees.
- Idea: divide (2D) space into four **quadrants**, and store items in the appropriate quadrant. Repeat this recursively with each quadrant that contains more than one item.
- Original definition: a quadtree is either
 - Empty, or
 - An item at some position (x, y) , called the root, plus
 - four quadtrees, each containing only items that are northwest, northeast, southwest, and southeast of (x, y) .
- Big idea is that if you are looking for point (x', y') and the root is not the point you are looking for, you can narrow down which of the four subtrees of the root to look in by comparing coordinates (x, y) with (x', y') .

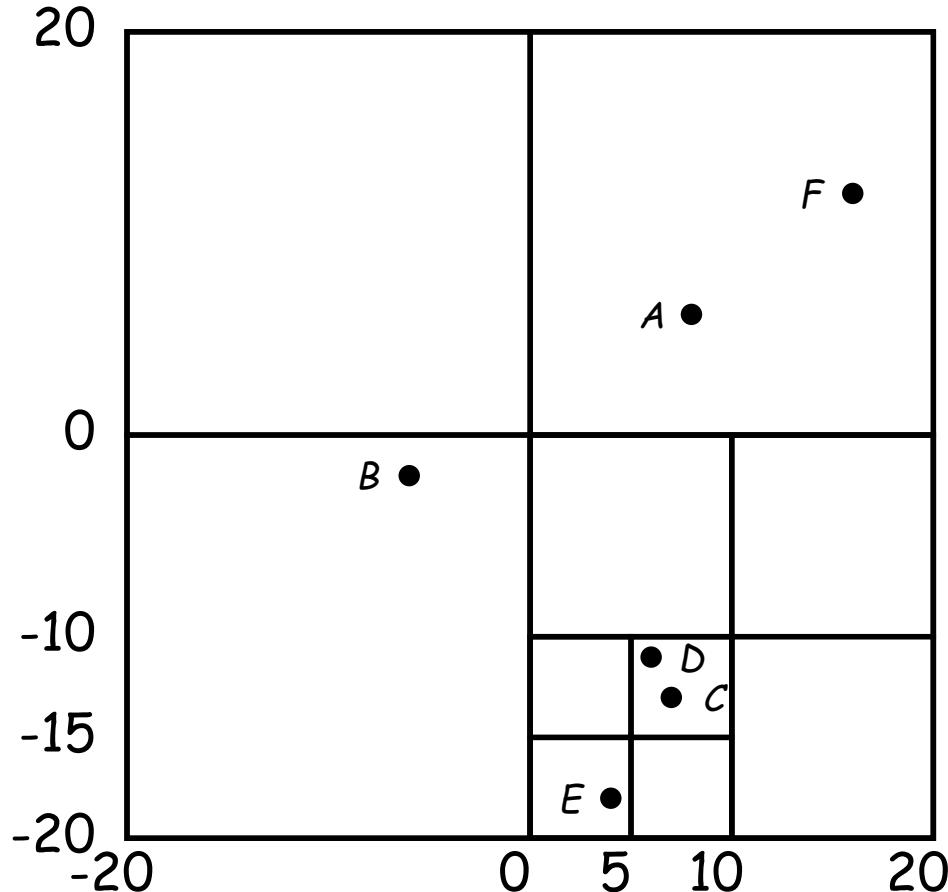
Classical Quadtree: Example



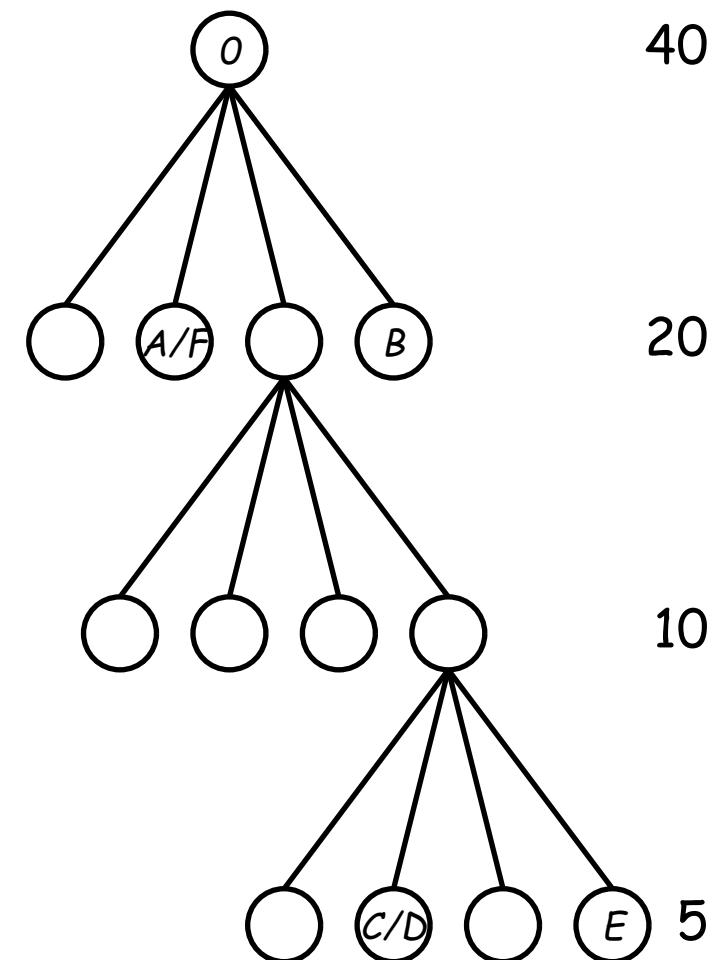
Point-region (PR) Quadtrees

- If we use a Quadtree to track moving objects, it may be useful to be able to *delete* items from a tree: when an object moves, the subtree that it goes in may change.
- Difficult to do with the classical data structure above, so we'll define instead:
- A quadtree consists of a bounding rectangle, B and either
 - Zero up to a small number of items that lie in that rectangle, or
 - Four quadtrees whose bounding rectangles are the four quadrants of B (all of equal size).
- A completely empty quadtree can have an arbitrary bounding rectangle, or you can wait for the first point to be inserted.

Example of PR Quadtree



(≤ 2 points per leaf)

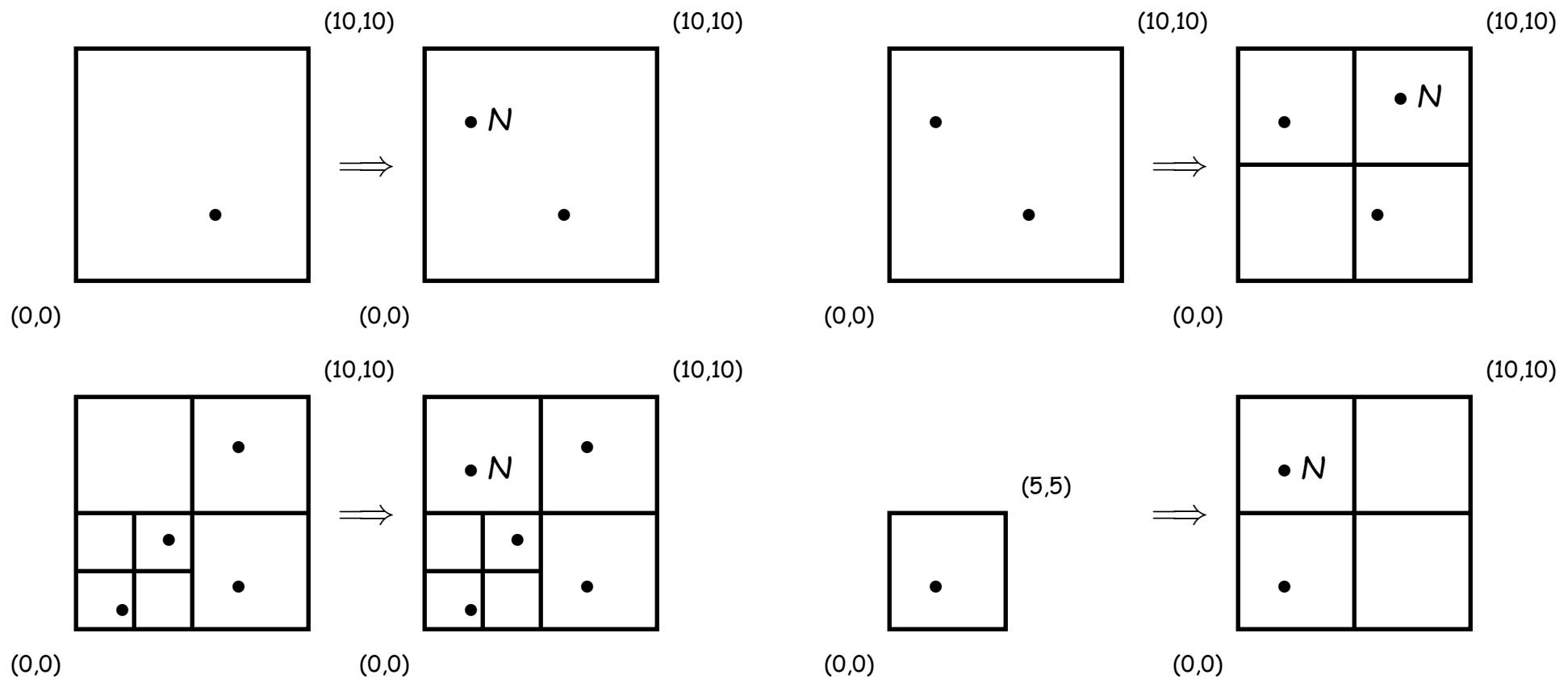


Navigating PR QuadTrees

- To find an item at (x, y) in quadtree T ,
 1. If (x, y) is outside the bounding rectangle of T , or T is empty, then (x, y) is not in T .
 2. Otherwise, if T contains a small set of items, then (x, y) is in T iff it is among these items.
 3. Otherwise, T consists of four quadtrees. Recursively look for (x, y) in each (however, step #1 above will cause all but one of these bounding boxes to reject the point immediately).
- Similar procedure works when looking for all items within some rectangle, R :
 1. If R does not intersect the bounding rectangle of T , or T is empty, then there are no items in R .
 2. Otherwise, if T contains a set of items, return those that are in R , if any.
 3. Otherwise, T consists of four quadtrees. Recursively look for points in R in each one of them.

Insertion into PR Quadtrees

Various cases for inserting a new point N , assuming maximum occupancy of a region is 2, showing initial state \Rightarrow final state.



CS61B Lecture #22: Hashing

Speeding Up Searching

- Linear search is OK for small data sets, bad for large.
- So linear search would be OK *if* we could rapidly narrow the search to a few items.
- Suppose that in constant time we could put any item in our data set into a numbered *bucket*, where # buckets stays within a constant factor of # keys.
- Suppose also that buckets contain roughly equal numbers of keys.
- Then search would be constant time in number of comparisons.

Hash functions

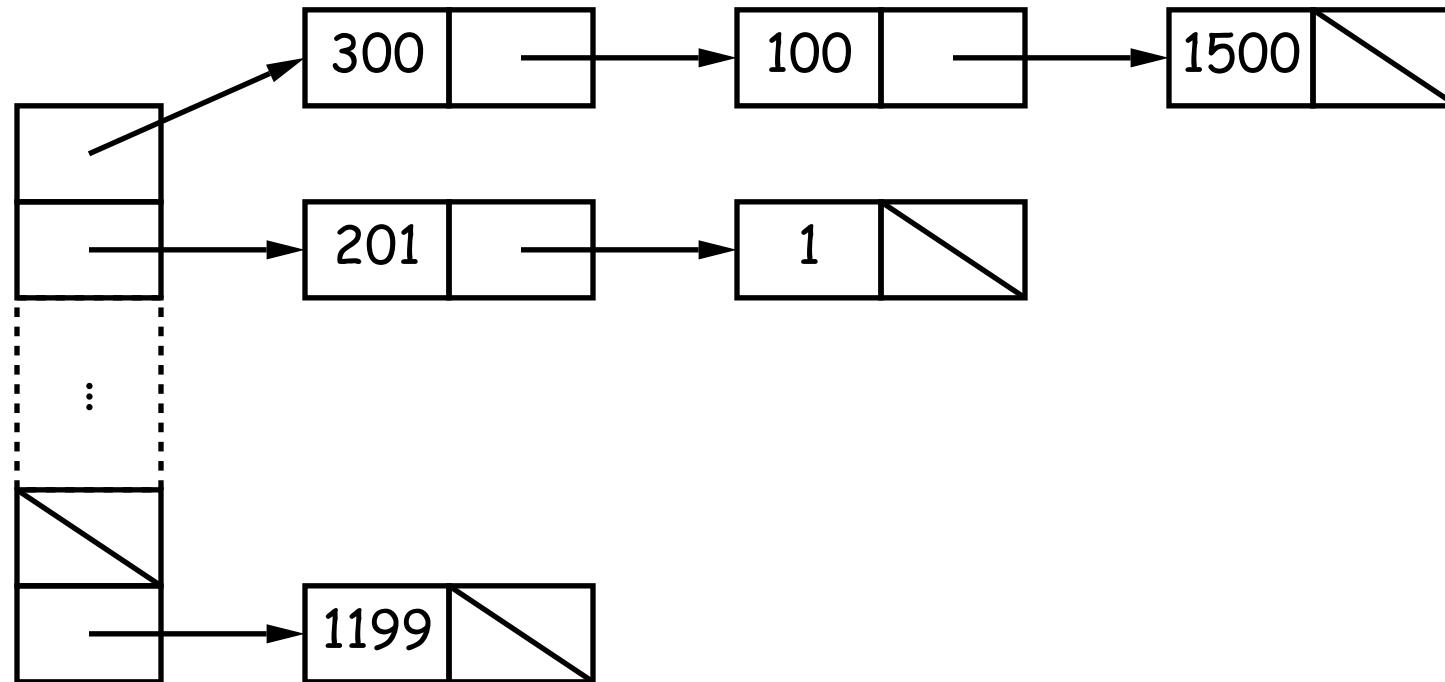
- To do this, must have way to convert key to bucket number: a *hash function*.

"hash /hæʃ/ 2a a mixture; a jumble. b a mess." *Concise Oxford Dictionary, eighth edition*

- Example:
 - $N = 200$ data items.
 - keys are longs, evenly spread over the range $0..2^{63} - 1$.
 - Want to keep maximum search to $L = 2$ items.
 - Use hash function $h(K) = K \% M$, where $M = N/L = 100$ is the number of buckets: $0 \leq h(K) < M$.
 - So 100232, 433, and 10002332482 go into different buckets, but 10, 400210, and 210 all go into the same bucket.

External chaining

- Array of M buckets.
- Each bucket is a list of data items.



- Not all buckets have same length, but average is $N/M = L$, the *load factor*.
- To work well, hash function must avoid too many *collisions*: keys that "hash" to equal values.

Ditching the Chains: Open Addressing

- Idea: Put one data item in each bucket.
- When there is a collision, and bucket is full, just use another.
- So you actually use a hash function with two arguments: $h(K, i)$, where K is the key, and the buckets you try are $h(K, 0), h(K, 1), \dots$
- Various possibilities (here S is the table size):
 - Linear probes: $h(K, i) = (h(K, 0) + C \cdot i) \bmod S$, for constant C .
 - Quadratic probes: $h(K, i) = (h(K, 0) + C_1 \cdot i + C_2 \cdot i^2) \bmod S$.
 - Double hashing: $h(K, i) = (h(K, 0) + i \cdot h'(K)) \bmod S$, where h' is a different hash function.
- Example: $h(K, i) = (K + i) \% M$, with $M = 10$:
 - Add 1, 2, 11, 3, 102, 9, 18, 108, 309 to empty table.

108	1	2	11	3	102	309		18	9
-----	---	---	----	---	-----	-----	--	----	---

- Things can get slow, even when table is far from full.
- Lots of literature on this technique, but
- Personally, I just settle for external chaining.

Filling the Table

- To get (likely to be) constant-time lookup, need to keep #buckets within constant factor of #items.
- So resize table when load factor gets higher than some limit.
- In general, must *re-hash* all table items.
- Still, this operation constant time per item,
- So by doubling table size each time, get constant *amortized* time for insertion and lookup
- (Assuming, that is, that our hash function is good).

Sets vs. Dictionaries

- So far, we've really been talking about hash tables containing *sets of keys*. The important operations are testing for membership, adding keys, and iterating through the keys.
- The term "hash table" can also refer to a *dictionary*, in which each key is associated with a *value*.
- The associated value does not affect in the search process.
- The important operations are testing for whether a key is in the table, finding what value is associated with a given key, and setting the value associated with a key.
- Here, we'll concentrate on sets of keys, since the addition of associated values is a conceptually small change.

Hash Functions: Strings

- For String, " $s_0s_1 \dots s_{n-1}$ " want function that takes all characters and their positions into account.
- What's wrong with $s_0 + s_1 + \dots + s_{n-1}$?
- For strings, Java uses

$$h(s) = s_0 \cdot 31^{n-1} + s_1 \cdot 31^{n-2} + \dots + s_{n-1}$$

computed modulo 2^{32} as in Java int arithmetic.

- To convert to a table index in $0..N - 1$, compute $h(s)\%N$ (but *don't* use table size that is multiple of 31!)
- Not as hard to compute as you might think; don't even need multiplication!

```
int r; r = 0;
for (int i = 0; i < s.length (); i += 1)
    r = (r << 5) - r + s.charAt (i);
```

Hash Functions: Other Data Structures I

- Lists (ArrayList, LinkedList, etc.) are analogous to strings: e.g., Java uses

```
hashCode = 1; Iterator i = list.iterator();
while (i.hasNext()) {
    Object obj = i.next();
    hashCode =
        31*hashCode
        + (obj==null ? 0 : obj.hashCode());
}
```

- Can limit time spent computing hash function by not looking at entire list. For example: look only at first few items (if dealing with a List or SortedSet).
- Causes more collisions, but does *not* cause equal things to go to different buckets.

Hash Functions: Other Data Structures II

- Recursively defined data structures \Rightarrow recursively defined hash functions.
- For example, on a binary tree, one can use something like

```
hash(T):  
    if (T == null)  
        return 0;  
    else return someHashFunction (T.label ())  
          ^ hash(T.left ()) ^ hash(T.right ());
```

Identity Hash Functions

- We can use the address of object as its hashed value ("hash on identity") if distinct (\neq) objects are never considered equal.
- But careful! Won't work for Strings, because .equal Strings could be in different buckets:

```
String H = "Hello",
      S1 = H + ", world!",
      S2 = "Hello, world!";
String H = "Hello",
      S1 = H + ", world!",
      S2 = "Hello, world!";
```

- Here `S1.equals(S2)`, but `S1 != S2`.
- Java provides the type `java.util.IdentityHashMap` for this purpose:

```
IdentityHashMap<Object, String> map = new IdentityHashMap<>();
Thing x = ...;
map.put(x, "Thing #1");    // Works if Thing does not override .equals.
map.put("Frank", "President"); // Probably not a good idea.
```

What Java Provides

- The class `Object` defines function `hashCode()`.
- By default, returns the identity hash function, or something similar.
[Why is this OK as a default?]
- Can override it for your particular type.
- For reasons given on last slide, it is overridden for type `String`, as well as many types in the Java library, like all kinds of `List`.
- The types `Hashtable`, `HashSet`, and `HashMap` use `hashCode` to give you fast look-up of objects.

```
HashMap<KeyType,ValueType> map =  
    new HashMap<>(approximate size, load factor);  
map.put(key, value);           // Map KEY -> VALUE.  
... map.get(someKey)           // VALUE last mapped to by SOMEKEY.  
... map.containsKey(someKey)   // Is SOMEKEY mapped?  
... map.keySet()              // All keys in MAP (a Set)
```

Special Case: Monotonic Hash Functions

- Suppose our hash function is *monotonic*: either nonincreasing or nondescreasing.
- So, e.g., if key $k_1 > k_2$, then $h(k_1) \geq h(k_2)$.
- Example:
 - Items are time-stamped records; key is the time.
 - Hashing function is to have one bucket for every hour.
- In this case, you *can* use a hash table to speed up range queries [How?]
- Could this be applied to strings? When would it work well?

Perfect Hashing

- Suppose set of keys is *fixed*.
- A tailor-made hash function might then hash every key to a different value: *perfect hashing*.
- In that case, there is no search along a chain or in an open-address table: either the element at the hash value is or is not equal to the target key.
- One technique: use a hash table of hash tables with two hash functions. One chooses a small hash table that uses the second hash function. With a few random selections of hash functions, you can soon find one where there is never a collision in the second hash table.
- For example, suppose we store key/value pairs in type Pair. Then:

```
Pair[][] table = ...;
// Find which table X must be in:
int h0 = hash0(X);
// Use different hash function for each h0.
Pair[] pair = table[h0][hash2(h0, X)];
return pair.key.equals(X) ? pair.value : null;
```

Characteristics

- Assuming good hash function, add, lookup, deletion take $\Theta(1)$ time, amortized.
- Good for cases where one looks up *equal* keys.
- For range queries such as "Give me every name between Martin and Napoli," hash tables are usually bad. [Why?]
- Hashing is probably not a good idea for small sets that you rapidly create and discard. [Why?]

CS61B Lecture #23

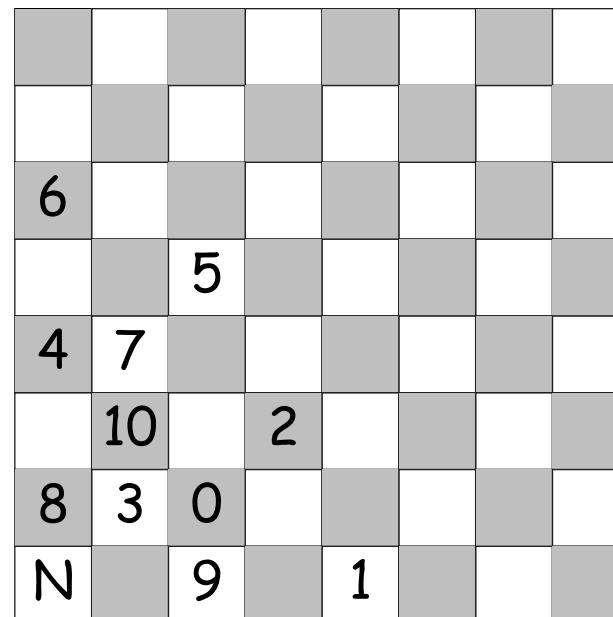
Today: Backtracking searches, game trees (*DSIJ*, Section 6.5)

Searching by “Generate and Test”

- We've been considering the problem of searching a set of data stored in some kind of data structure: "Is $x \in S$?"
- But suppose we *don't* have a set S , but know how to recognize what we're after if we find it: "Is there an x such that $P(x)$?"
- If we know how to enumerate all possible candidates, can use approach of *Generate and Test*: test all possibilities in turn.
- Can sometimes be more clever: avoid trying things that won't work, for example.
- What happens if the set of possible candidates is infinite?

Backtracking Search

- Backtracking search is one way to enumerate all possibilities.
- Example: *Knight's Tour*. Find all paths a knight can travel on a chessboard such that it touches every square exactly once and ends up one knight move from where it started.
- In the example below, the numbers indicate position numbers (knight starts at 0).
- Here, knight (N) is stuck; how to handle this?



General Recursive Algorithm

```
/** Append to PATH a sequence of knight moves starting at ROW, COL
 * that avoids all squares that have been hit already and
 * that ends up one square away from ENDROW, ENDCOL. B[i][j] is
 * true iff row i and column j have been hit on PATH so far.
 * Returns true if it succeeds, else false (with no change to PATH
 * or B). Call initially with PATH containing the starting square, and
 * the starting square (only) marked in B. */

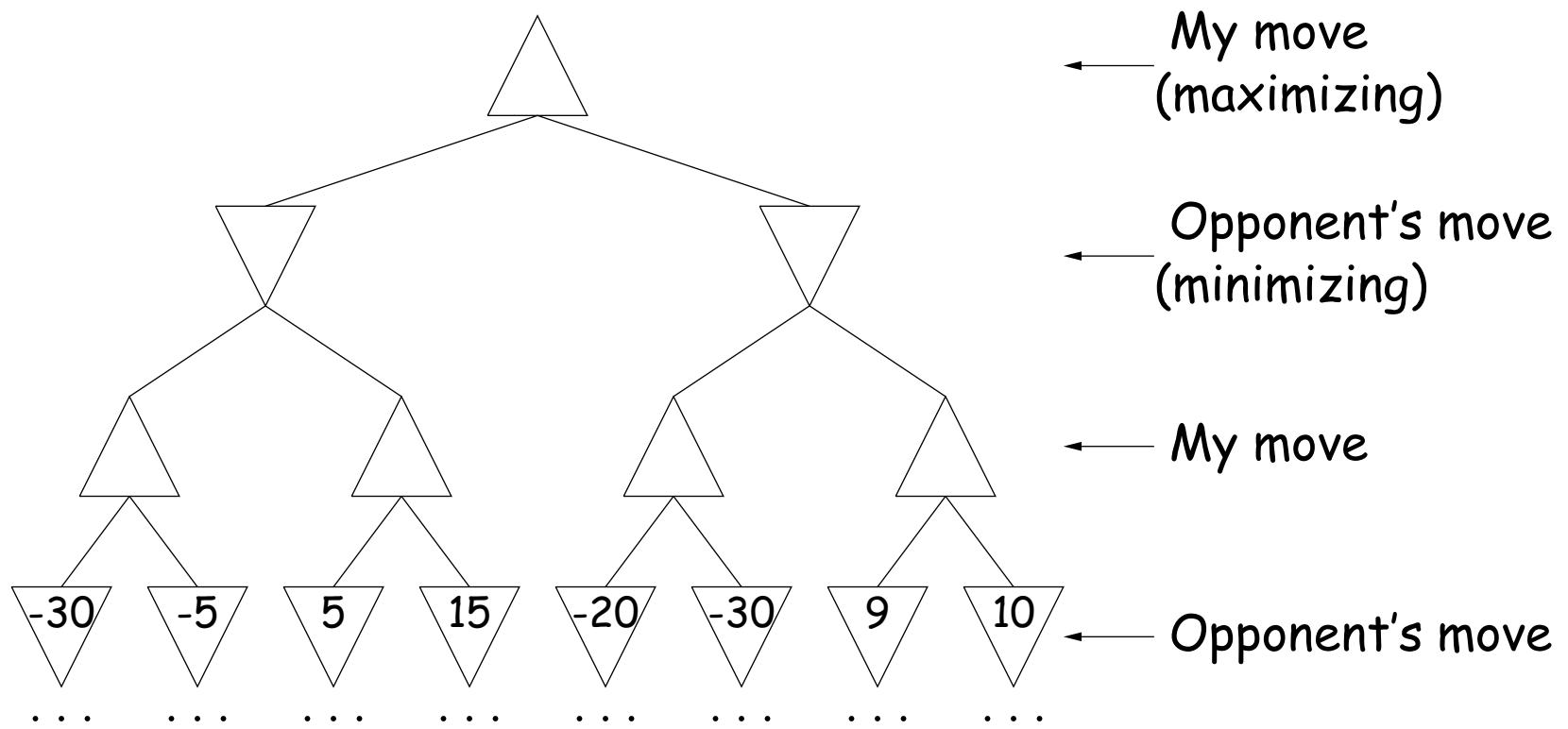
boolean findPath(boolean[][] b, int row, int col,
                 int endRow, int endCol, List path) {
    if (path.size() == 64)    return isKnightMove(row, col, endRow, endCol);
    for (r, c = all possible moves from (row, col)) {
        if (!b[r][c]) {
            b[r][c] = true; // Mark the square
            path.add(new Move(r, c));
            if (findPath(b, r, c, endRow, endCol, path))  return true;
            b[r][c] = false; // Backtrack out of the move.
            path.remove(path.size()-1);
        }
    }
    return false;
}
```

Another Kind of Search: Best Move

- Consider the problem of finding the *best* move in a two-person game.
- One way: assign a *heuristic value* to each possible move and pick highest (aka *static valuation*).
- Otherwise, we can use a variety of heuristics. Some examples of static valuations:
 - assign a maximal or minimal value to a won position (depending on side.)
 - number of black pieces – number of white pieces in checkers.
 - (weighted sum of white piece values) – (weighted sum of black pieces in chess), such as queen=9, rook=5, knight=bishop=3, pawn=1.
 - Nearness of pieces to strategic areas (center of board).
- But this is misleading. A move might give us more pieces, but set up a devastating response from the opponent.
- So, for each move, look at *opponent's* possible moves, use the best move that results *for the opponent* as the value.
- But what if *you* have a great response to opponent's response?
- How do we organize this sensibly?

Game Trees

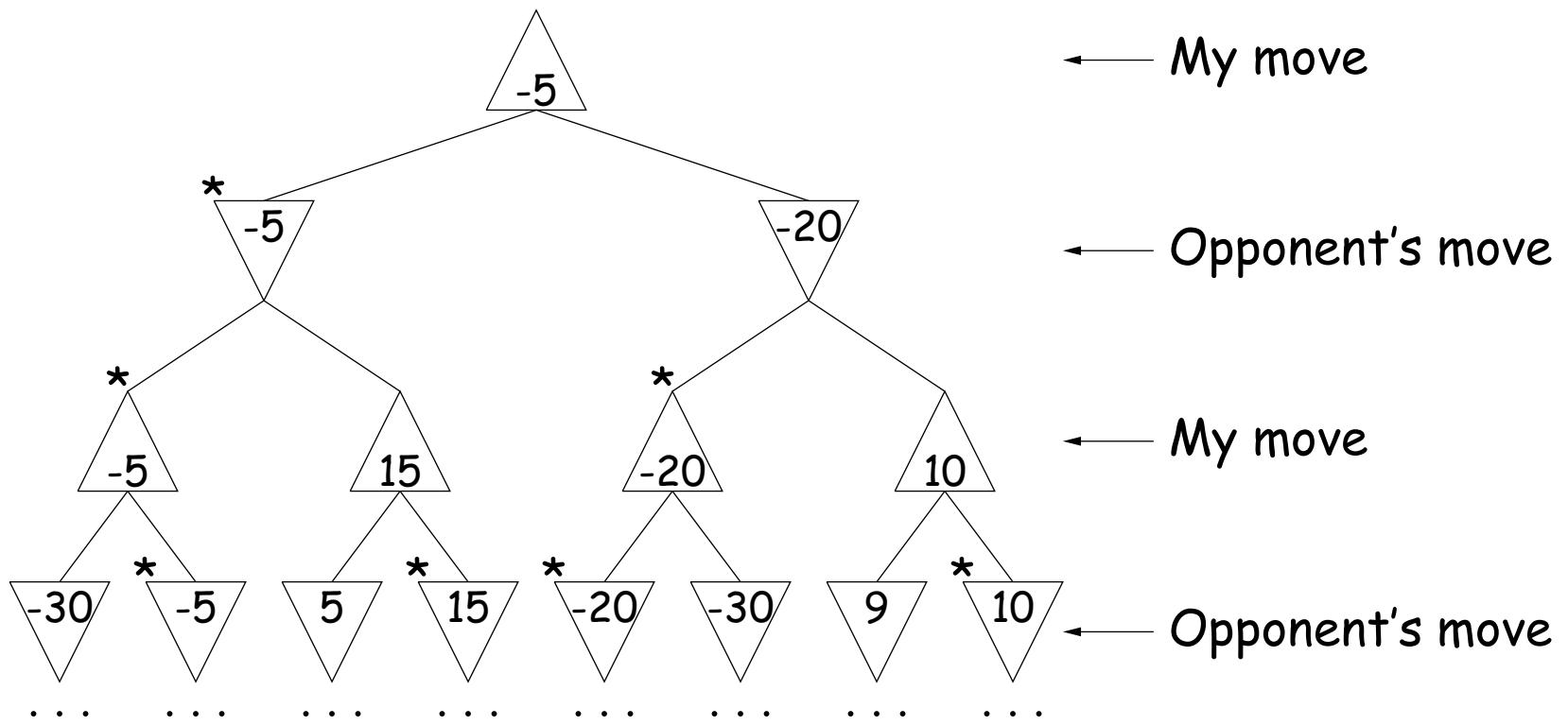
- Think of the space of possible continuations of the game as a tree.
- Each node is a position, each edge a move. Suppose numbers at the bottom are the values of those positions *to me*. Smaller numbers are of more value to *my opponent*.



- What value can I get if my opponent plays as well as possible? How?

Game Trees, Minimax

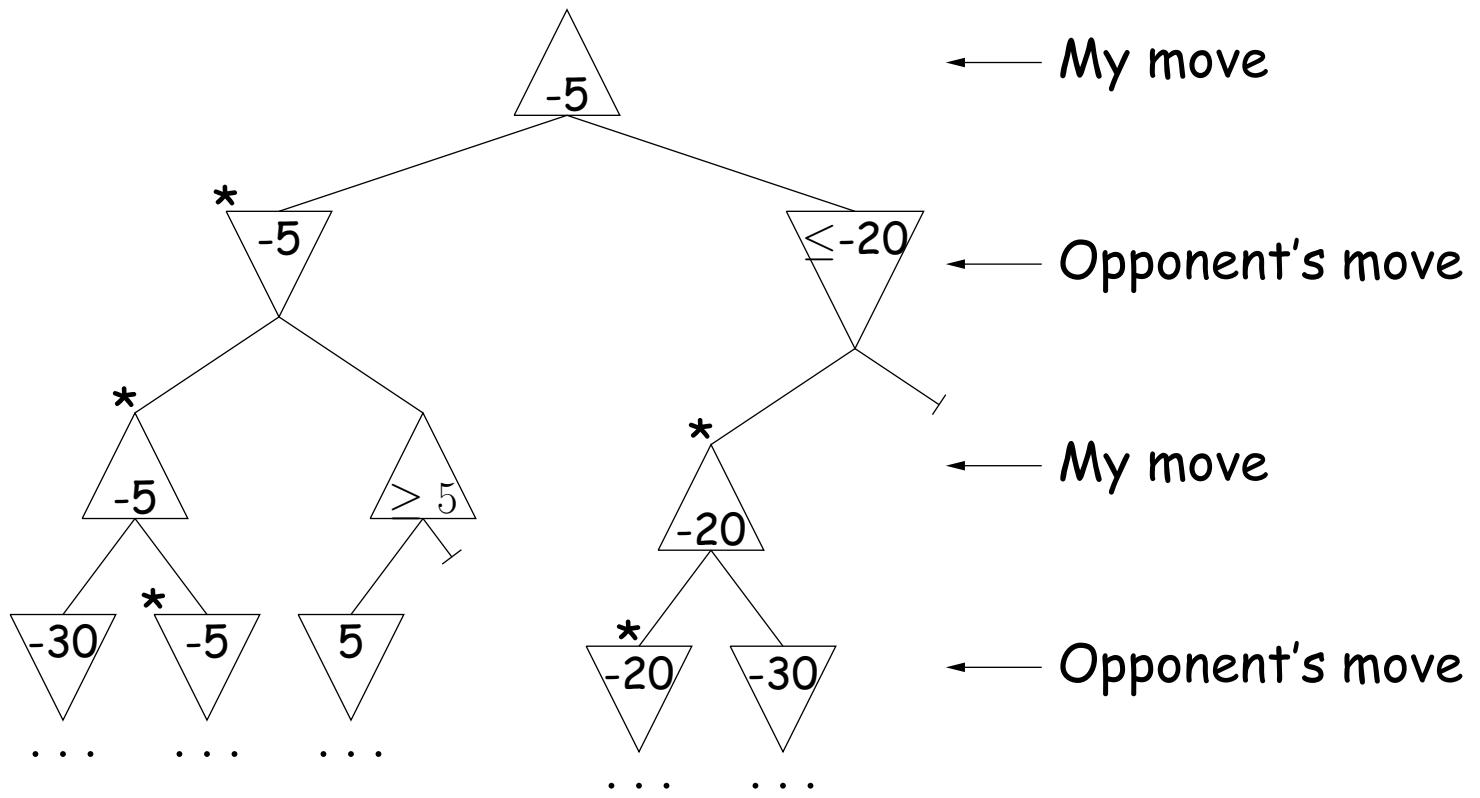
- Think of the space of possible continuations of the game as a tree.
- Each node is a position, each edge a move. Numbers are the values we guess for the positions (larger means better for me). Starred nodes would be chosen.



- I always choose child (next position) with maximum value; opponent chooses minimum value—the **minimax algorithm**.

Alpha-Beta Pruning

- We can *prune* this tree as we search it.



- At the ' ≥ 5 ' position, I know that the opponent will not choose to move here (already has a -5 move).
- At the ' ≤ -20 ' position, my opponent knows that I will never choose to move here (since I already have a -5 move).

Cutting off the Search

- If you could traverse game tree to the bottom, you'd be able to force a win (if it's possible).
- Sometimes possible near the end of a game.
- Unfortunately, game trees tend to be either infinite or impossibly large.
- So, we choose a maximum *depth*, and use a heuristic static valuation as the value at that depth.
- Or we might use *iterative deepening*, repeating the search at increasing depths until time is up.
- Much more sophisticated searches are possible, however (take CS188).

Overall Search Algorithm

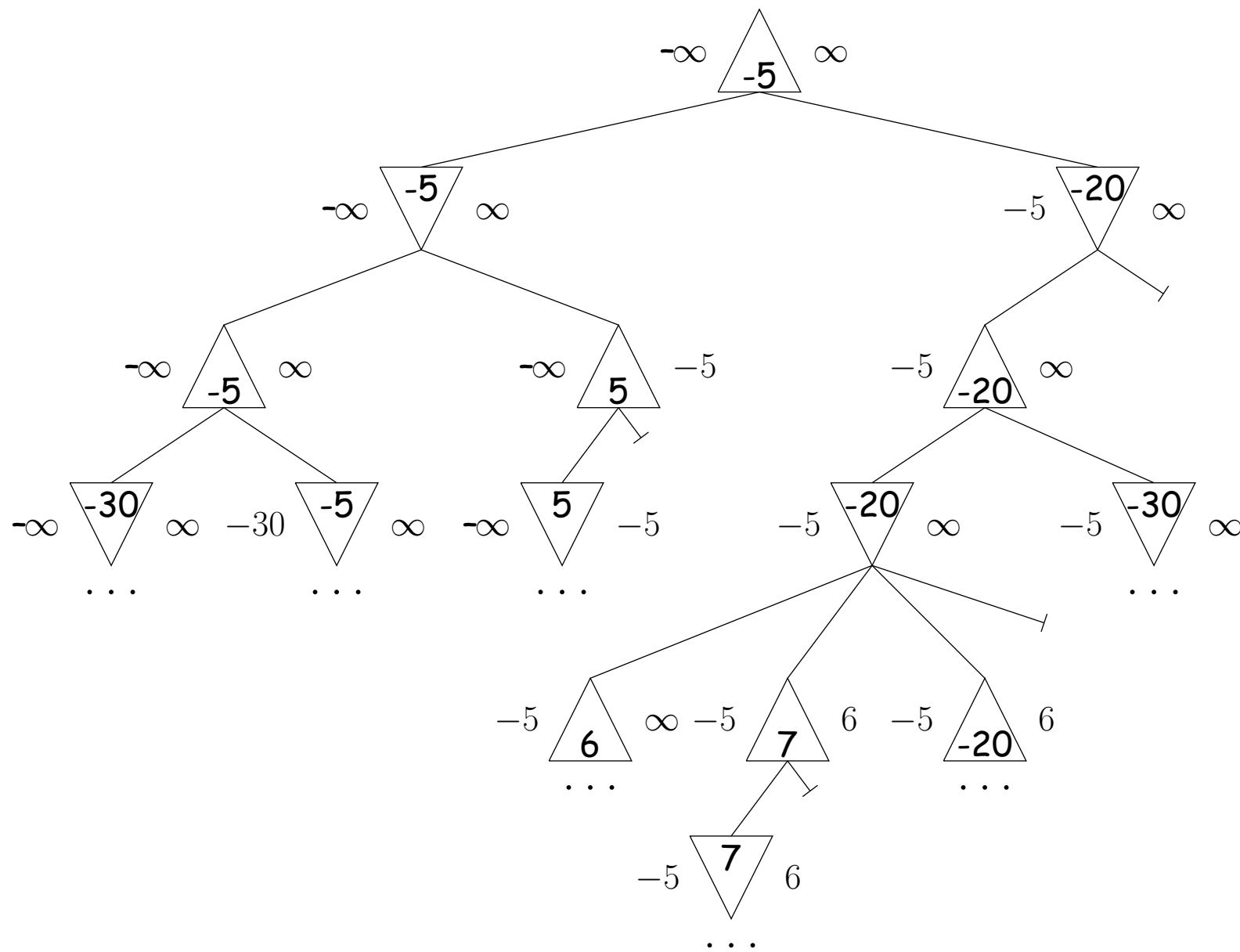
- Depending on whose move it is (maximizing player or minimizing player), we'll search for a move estimated to be optimal in one direction or the other.
- Search will be exhaustive down to a particular depth in the game tree; below that, we guess values.
- Also pass α and β limits:
 - High player does not care about exploring a position further after finding that its value will be larger than a position the minimizing player has already found, because the minimizing player will simply not choose a position with that larger value.
 - Likewise, minimizing player won't explore a positions whose value is less than what the maximizing player can get (α).
- To start, a maximizing player will find a move with the call

`maxPlayerValue(current position, search depth, $-\infty, +\infty$)`

- minimizing player:

`minPlayerValue(current position, search depth, $-\infty, +\infty$)`

Sample Tree with Alpha and Beta Values



Some Pseudocode for Searching (Maximizing Player)

```
/** The estimated minimax value of position POSN, searching up to
 * DEPTH moves ahead, assuming it is the maximizing player's move.
 * If the value is determined to be <=ALPHA, then the function
 * may return any value <=ALPHA, even if inaccurate. Likewise if the
 * value is >=BETA, it may return any value >=BETA. Assumes ALPHA<BETA. */
int maxPlayerValue(Position posn, int depth, int alpha, int beta)
{
    if (posn is a final position of the game || depth == 0)
        return staticGuess(posn);
    int bestSoFar = -∞;
    for (each legal move, M, in position posn) {
        Position next = makeMove(posn, M);
        int response = minPlayerValue(next, depth-1, alpha, beta);
        if (response > bestSoFar) {
            bestSoFar = response;
            alpha = max(alpha, bestSoFar);
            if (alpha >= beta)
                return bestSoFar;
        }
    }
    return bestSoFar;
}
```

Some Pseudocode for Searching (Minimizing Player)

```
/** The estimated minimax value of position POSN, searching up to
 * DEPTH moves ahead, assuming it is the minimizing player's move. */
int minPlayerValue(Position posn, int depth, int alpha, int beta)
{
    if (posn is a final position of the game || depth == 0)
        return staticGuess(posn);
    int bestSoFar = +∞;
    for (each legal move, M, in position posn) {
        Position next = makeMove(posn, M);
        int response = maxPlayerValue(next, depth-1, alpha, beta);
        if (response < bestSoFar) {
            bestSoFar = response;
            beta = min(beta, bestSoFar);
            if (alpha >= beta)
                return bestSoFar;
        }
    }
    return bestSoFar;
}
```

CS61B Lecture #24

Today:

- Priority queues (*Data Structures* §6.4, §6.5)
- Range queries (§6.2)
- Java utilities: SortedSet, Map, etc.

Priority Queues, Heaps

- Priority queue: defined by operations "add," "find largest," "remove largest."
- Examples: scheduling long streams of actions to occur at various future times.
- Also useful for sorting (keep removing largest).
- Common implementation is the **heap**, a kind of tree.
- (Confusingly, this same term is used to described the pool of storage that the **new** operator uses. Sorry about that.)

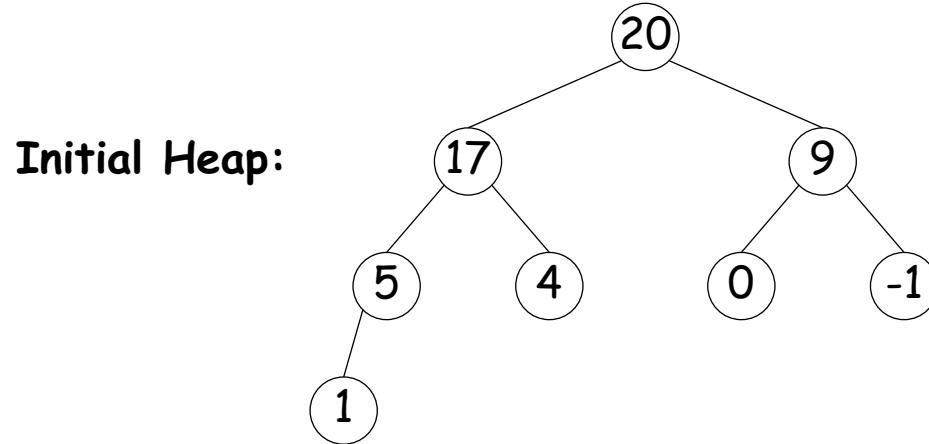
Heaps

- A *max-heap* is a binary tree that enforces the *Heap Property*: Labels of *both* children of each node are less than node's label.
- So the node at top has largest label.
- Looser than binary search property, which allows us to keep tree "bushy".
- That is, it's always valid to put the smallest nodes anywhere at the bottom of the tree.
- Thus, heaps can be made *nearly complete*: all but possibly the last row have as many keys as possible.
- As a result, insertion of new value and deletion of largest value always take time proportional to $\lg N$ in worst case.
- A *min-heap* is basically the same, but with the minimum value at the root and children having larger values than their parents.

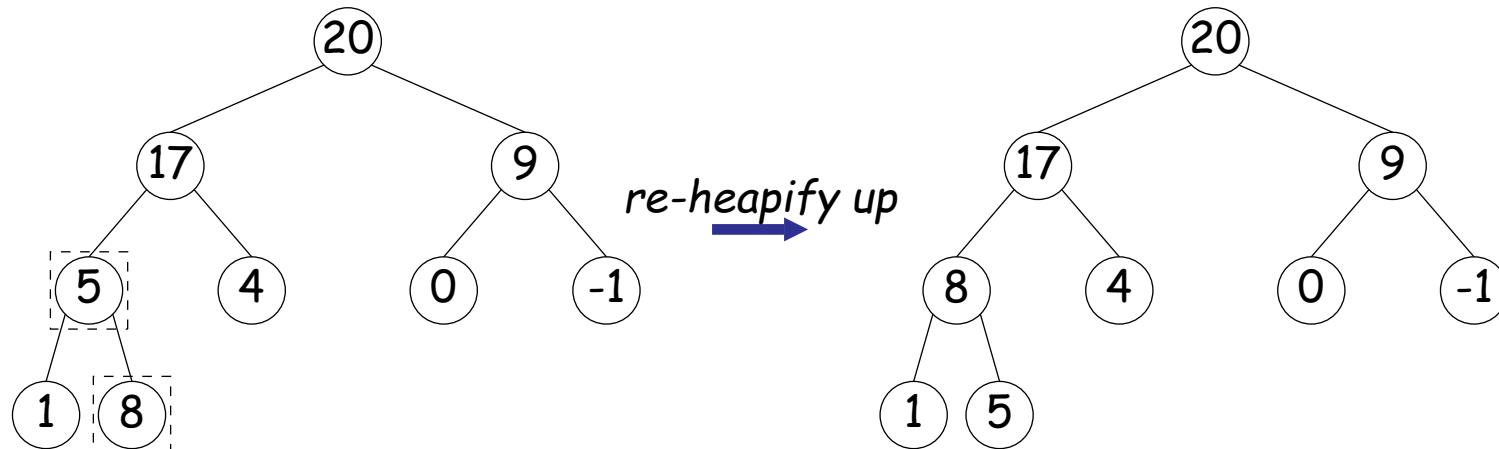
Example: Inserting into a simple heap

Data:

1 17 4 5 9 0 -1 20

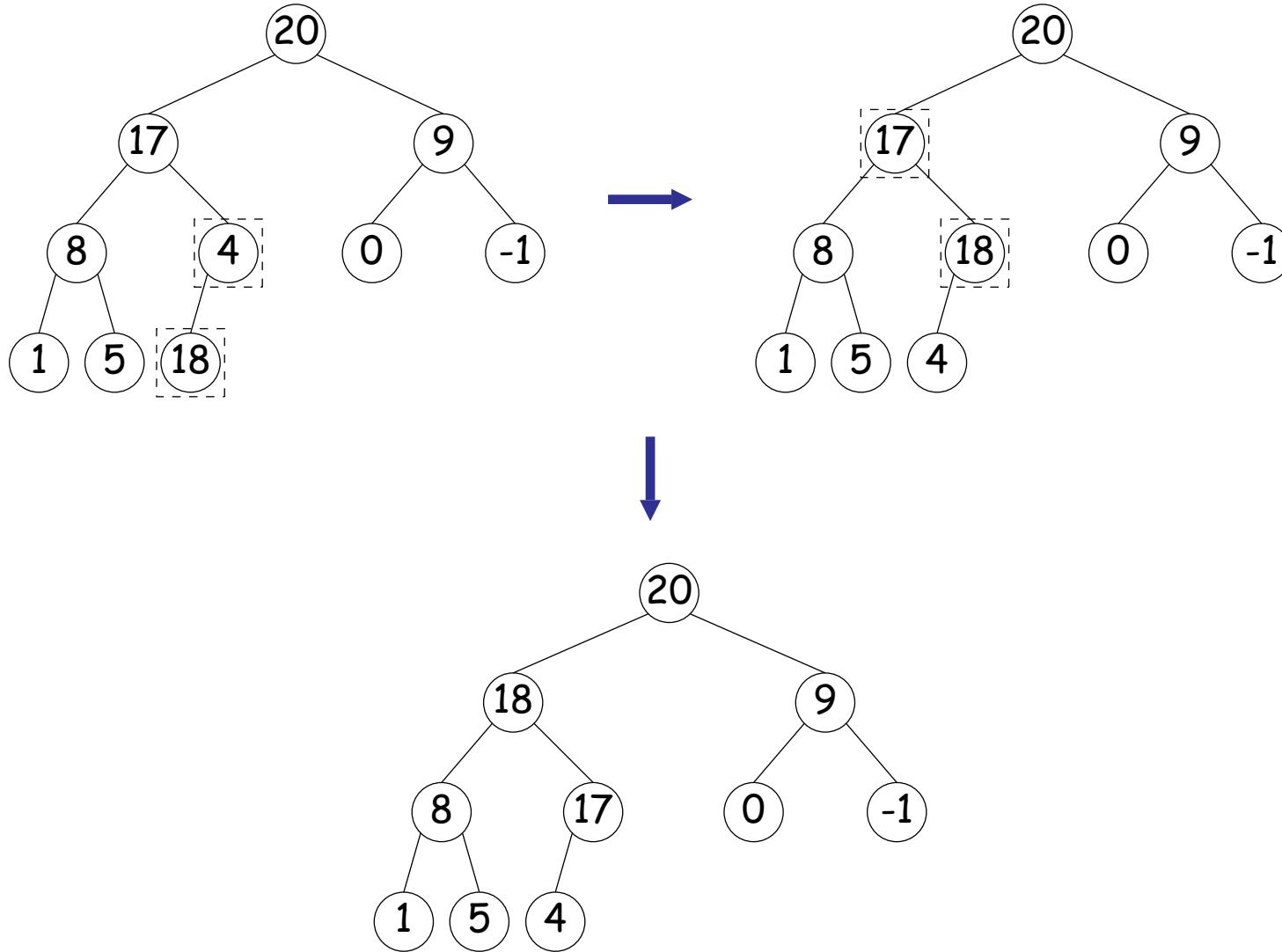


Add 8: Dashed boxes show where the heap property is violated



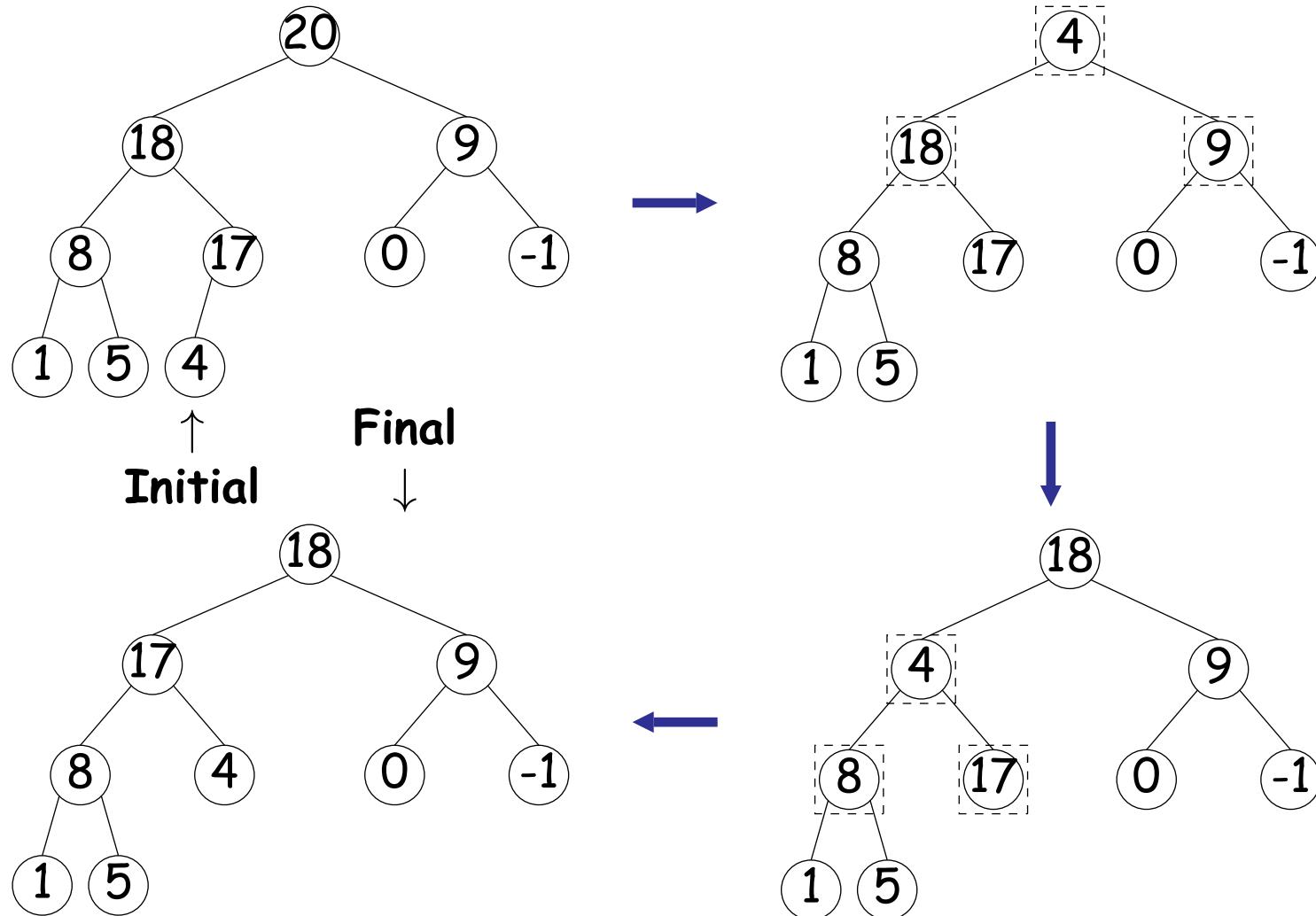
Heap insertion continued

Now insert 18:



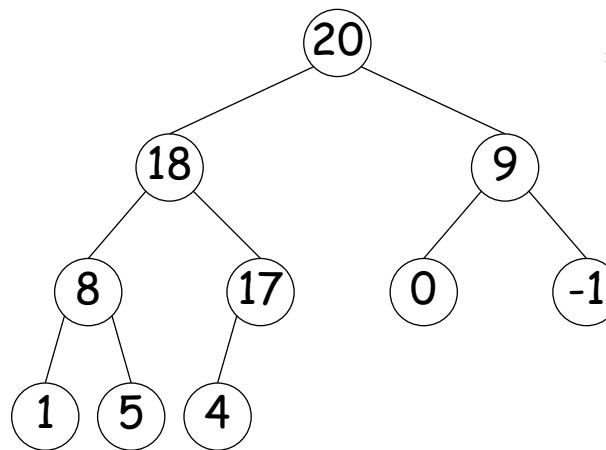
Removing Largest from Heap

To remove largest: Move the bottommost, rightmost node to the top, then re-heapify down as needed (swap offending node with larger child) to re-establish the heap property.



Heaps in Arrays

- Since heaps are nearly complete (missing items only at bottom level), can use arrays for compact representation.
- Example of removal from last slide (dashed arrows show children):



⇒

1	2	3	4	5	6	7	8	9	10
20	18	9	8	17	0	-1	1	5	4



4	18	9	8	17	0	-1	1	5	
---	----	---	---	----	---	----	---	---	--



18	4	9	8	17	0	-1	1	5	
----	---	---	---	----	---	----	---	---	--



18	17	9	8	4	0	-1	1	5	
----	----	---	---	---	---	----	---	---	--

Nodes stored in level order.

Children of node at index # K are in $2K$ and $2K + 1$ if numbering from 1, or $2K + 1$ and $2K + 2$ if from 0.

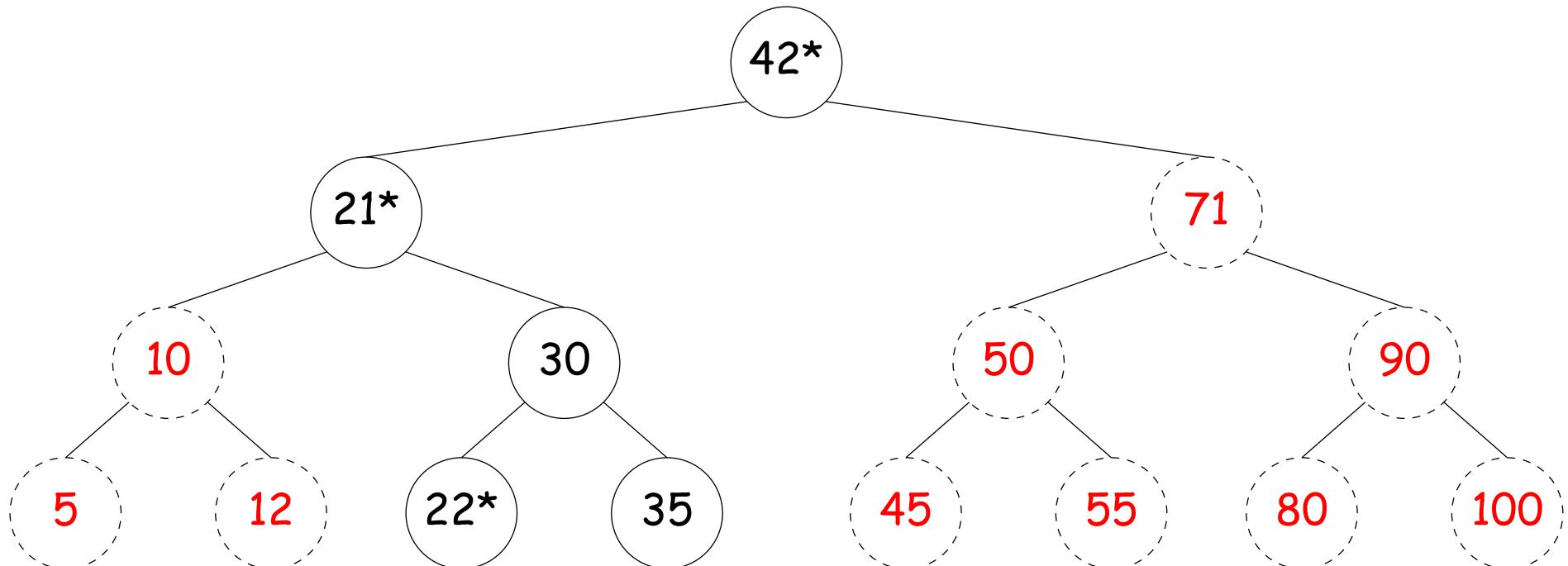
Ranges

- So far, we have looked for specific items,
- But for BSTs, we need an ordering anyway, and we can also support looking for *ranges of values*.
- Example: perform some action on all values in a BST that are within some range (in natural order):

```
/** Apply WHATTODO to all labels in T that are >= L and < U,
 * in ascending natural order. */
static void visitRange(BST<String> T, String L, String U,
                      Consumer<BST<String>> whatToDo) {
    if (T != null) {
        int compLeft = L.compareTo(T.label()),
            compRight = U.compareTo(T.label());
        if (compLeft < 0)                  /* L < label */
            visitRange (T.left(), L, U, whatToDo);
        if (compLeft <= 0 && compRight > 0) /* L <= label < U */
            whatToDo.accept(T);
        if (compRight > 0)                /* label < U */
            visitRange (T.right (), L, U, whatToDo);
    }
}
```

Time for Range Queries

- Time for range query $\in O(h + M)$, where h is the height of the tree, and M is the number of data items that turn out to be in the range.
- Consider searching the tree below for all values $25 \leq x < 40$.
- **Dashed** nodes are never looked at. Starred nodes are looked at but not output. The h comes from the starred nodes; the M comes from unstarred non-dashed nodes.



Ordered Sets and Range Queries in Java

- Class SortedSet supports range queries with *views* of a set:
 - S.headSet(U): subset of S that is $< U$.
 - S.tailSet(L): subset that is $\geq L$.
 - S.subSet(L, U): subset that is $\geq L, < U$.
- Changes to views modify S.
- Attempts to, e.g., add to a headSet beyond U are disallowed.
- Can iterate through a view to process a range:

```
SortedSet<String> fauna = new TreeSet<String>
    (Arrays.asList ("axolotl", "elk", "dog", "hartebeest", "duck"));
for (String item : fauna.subSet ("bison", "gnu"))
    System.out.printf ("%s, ", item);
```

would print "dog, duck, elk,"

TreeSet

- Java library type `TreeSet<T>` requires either that `T` be Comparable, or that you provide a Comparator, as in:

```
SortedSet<String> rev_fauna = new TreeSet<String>(Collections.reverseOrder());
```

- Comparator is a type of function object:

```
interface Comparator<T> {  
    /** Return <0 if LEFT<RIGHT, >0 if LEFT>RIGHT, else 0. */  
    int compare(T left, T right);  
}
```

(We'll deal with what `Comparator<T>` extends `Comparable<T>` is all about later.)

- For example, the `reverseOrder` comparator is defined like this:

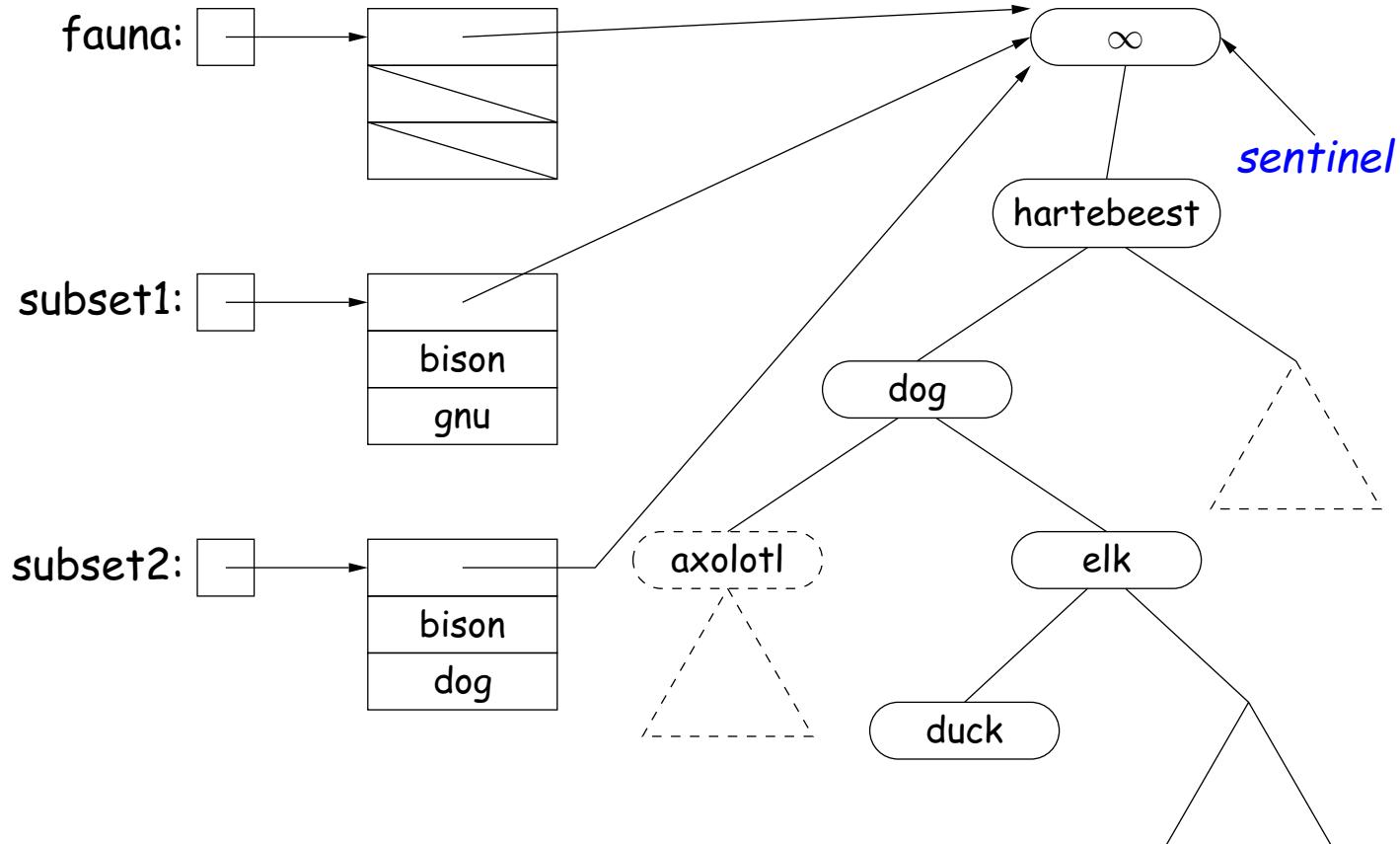
```
/** A Comparator that gives the reverse of natural order. */  
static <T extends Comparable<T>> Comparator<T> reverseOrder() {  
    // Java figures out this lambda expression is a Comparator<T>.   
    return (x, y) -> y.compareTo(x);  
}
```

Example of Representation: BSTSet

- Same representation for both sets and subsets.
- Pointer to BST, plus bounds (if any).
- `.size()` is expensive!

SortedSet<String>

```
fauna = new BSTSet<String>(stuff);  
subset1 = fauna.subSet("bison", "gnu");  
subset2 = subset1.subSet("axolotl", "dog");
```



Comparing Search Structures

Here, N is #items, k is #answers to query.

Function	Unordered	Sorted	Bushy Search	"Good" Hash Table	Heap
	List	Array	Tree		
<i>find</i>	$\Theta(N)$	$\Theta(\lg N)$	$\Theta(\lg N)$	$\Theta(1)$	$\Theta(N)$
<i>add (amortized)</i>	$\Theta(1)$	$\Theta(N)$	$\Theta(\lg N)$	$\Theta(1)$	$\Theta(\lg N)$
<i>range query</i>	$\Theta(N)$	$\Theta(k + \lg N)$	$\Theta(k + \lg N)$	$\Theta(N)$	$\Theta(N)$
<i>find largest</i>	$\Theta(N)$	$\Theta(1)$	$\Theta(\lg N)$	$\Theta(N)$	$\Theta(1)$
<i>remove largest</i>	$\Theta(N)$	$\Theta(1)$	$\Theta(\lg N)$	$\Theta(N)$	$\Theta(\lg N)$

CS61B Lecture #25: Java Generics

The Old Days

- Java library types such as List didn't used to be parameterized. All Lists were lists of Objects.
- So you'd write things like this:

```
void max(List L) {  
    String s; s = "";  
    for (int i = 0; i < L.size(); i += 1)  
        if (s.compareTo((String) L.get(i)) < 0) s = (String) L.get(i);  
    return s;  
}
```

- That is, you had to explicitly cast result of L.get(i) to let the compiler know what its static type was.
- Also, when calling L.add(x), there was no check that you put only Strings into it.
- So, starting with 1.5, the designers tried to alleviate these perceived problems by introducing **parameterized types**, like List<String>.
- Unfortunately, it is not as simple as one might think.

Basic Parameterization

- From the definitions of ArrayList and Map in java.util:

```
public class ArrayList<Item> implements List<Item> {  
    public Item get(int i) { ... }  
    public boolean add(Item x) { ... }  
    ...  
}  
public interface Map<Key, Value> {  
    Value get(Key x);  
    ...  
}
```

- First (boxed) occurrences of Item, Key, and Value introduce formal *type parameters*, whose “values” (which are reference types) get substituted for all the other occurrences of Item, Key, or Value when ArrayList or Map is “called” (as in ArrayList<String>, or ArrayList<int []>, or Map<String, List<Particle>>).
- Other occurrences of Item, Key, and Value are *uses* of the formal types, just like uses of a formal parameter in the body of a function.

Type Instantiation

- *Instantiating* a generic type is analogous to calling a function.
- Consider again

```
public class ArrayList<Item> implements List<Item> {  
    public Item get(int i) { ... }  
    public boolean add(Item x) { ... } ...  
}
```

- When we write `ArrayList<String>`, we get, in effect, a new type and a new interface, somewhat like

```
public interface String_List {  
    String get(int i);  
    boolean add(String x); ...  
}
```

```
public String_ArrayList implements String_List {  
    public String get(int i) { ... }  
    public boolean add(String x) { ... } ...  
}
```

Parameters on Methods

- Functions (methods) may also be parameterized by type. Example of use from `java.util.Collections`:

```
/** A read-only list containing just ITEM. */
static <T> List<T> singleton(T item) { ... }
/** An unmodifiable empty list. */
static <T> List<T> emptyList() { ... }
```

The compiler figures out T in the expression `singletont(x)` by looking at the type of `x`. This is a simple example of *type inference*.

- In the call

```
List<String> empty = Collections.emptyList();
```

the parameters obviously don't suffice, but the compiler deduces the parameter `T` from context: it must be assignable to `String`.

Wildcards

- Consider the definition of something that counts the number of times something occurs in a collection of items. We could write:

```
/** Number of items in C that are equal to X. */
static <T> int frequency(Collection<T> c, Object x) {
    int n; n = 0;
    for (T y : c) {
        if (x.equals(y))
            n += 1;
    }
    return n;
}
```

- But we don't really care what T is; we don't need to declare anything of type T in the body, because we could write instead

```
...
for (Object y : c) {
```

- *Wildcard type parameters* say that you don't care what a type parameter is (i.e., it's any subtype of Object):

```
static int frequency(Collection<?> c, Object x) {...}
```

Subtyping (I)

- What are the relationships between the types

`List<String>, List<Object>, ArrayList<String>, ArrayList<Object>?`

- We know that `ArrayList` \preceq `List` and `String` \preceq `Object` (using \preceq for "is a subtype of")...
- ... So is `List<String>` \preceq `List<Object>?`

Subtyping (II)

- Consider this fragment:

```
List<String> LS = new ArrayList<String>();  
List<Object> LObj = LS;      // OK??  
int[] A = { 1, 2 };  
LObj.add(A);                // Legal, since A is an Object  
String S = LS.get(0);        // OOPS! A.get(0) is NOT a String,  
                           // but spec of List<String>.get  
                           // says that it is.
```

- So, having $\text{List} < \text{String} \preceq \text{List} < \text{Object}$ would violate *type safety*:
The compiler is wrong about the type of a value.
- So in general for $T_1 < X \preceq T_2 < Y$, must have $X = Y$.
- But what about T_1 and T_2 ?

Subtyping (III)

- Now consider

```
ArrayList<String> ALS = new ArrayList<String>();  
List<String> LS = ALS;           // OK??
```

- In this case, everything's fine:

- The object's dynamic type is `ArrayList<String>`.
- Therefore, the methods expected for `LS` must be a subset of those for `ALS`.
- And since the type parameters are the same, the signatures of those methods will be the same.
- Therefore, all the legal calls on methods of `LS` (according to the compiler) will be valid for the actual object pointed to by `LS`.

- In general, $T1<X> \preceq T2<X>$ if $T1 \preceq T2$.

A Java Inconsistency: Arrays

- The Java language design is not entirely consistent when it comes to subtyping.
- For the same reason that `ArrayList<String>` $\not\subseteq$ `ArrayList<Object>`, you'd also expect that `String[]` $\not\subseteq$ `Object[]`.
- And yet, Java *does* make `String[]` \subseteq `Object[]`.
- And, just as explained above, one gets into trouble with

```
String[] AS = new String[3];
Object[] AObj = AS;
AObj[0] = new int[] { 1, 2 }; // Bad
```

- So in Java, the **Bad** line causes an `ArrayStoreException`—a (dynamic) runtime error instead of a (static) compile-time error.
- Why do it this way? Basically, because otherwise there'd be no way to implement, e.g., `ArrayList`. (For reasons to be explained, you can't use the type parameter to `ArrayList` as the type of array elements, nor as the type of a cast.)

Type Bounds (I)

- Sometimes, your program needs to ensure that a particular type parameter is replaced only by a subtype (or supertype) of a particular type (sort of like specifying the "type of a type.").
- For example,

```
class NumericSet<T extends Number> extends HashSet<T> {  
    /** My minimal element */  
    T min() { ... }  
    ...  
}
```

Requires that all type parameters to NumericSet must be subtypes of Number (the "type bound"). T can either extend or implement the bound, as appropriate.

Type Bounds (II)

- Another example:

```
/** Copy the elements of SRC into DEST. */
static <T> void copy(List<? super T> dest, List<T> src) { ... }
```

means that dest can be a List<Q> for any Q as long as T is a subtype of (extends or implements) Q.

- Why not just define this as

```
/** Copy the elements of SRC into DEST. */
static <T> void copy(List<T> dest, List<T> src) { ... }
```

?

Type Bounds (II)

- Another example:

```
/** Copy the elements of SRC into DEST. */
static <T> void copy(List<? super T> dest, List<T> src) { ... }
```

means that dest can be a List<Q> for any Q as long as T is a subtype of (extends or implements) Q.

- Why not just define this as

```
/** Copy the elements of SRC into DEST. */
static <T> void copy(List<T> dest, List<T> src) { ... }
```

?

- It would make perfect sense to copy a List<String> into a List<Object>, but that would be disallowed with this declaration.

Type Bounds (III)

- And one more:

```
/** Search sorted list L for KEY, returning either its position (if
 *  present), or k-1, where k is where KEY should be inserted. */
static <T> int binarySearch(List<? extends Comparable<? super T>> L,
                           T key)
```

- Here, the items of L have to have a type that is comparable to T's or to some supertype of T.
- Does L have to be able to contain the value key?
- Why does this make sense?

Type Bounds (III)

- And one more:

```
/** Search sorted list L for KEY, returning either its position (if
 *  present), or k-1, where k is where KEY should be inserted. */
static <T> int binarySearch(List<? extends Comparable<? super T>> L,
                           T key)
```

- Here, the items of L have to have a type that is comparable to T's or to some supertype of T.
- Does L have to be able to contain the value key?
- Why does this make sense?
- As long as the items in L can be compared to key, it doesn't really matter whether they might include key (not that this is often useful).

Dirty Secrets Behind the Scenes

- Java's design for parameterized types was constrained by a desire for backward compatibility.
- They wanted to be able to use the same interpreter as before for programs with these new generic types, so it had to be possible to rewrite generic programs into old-style non-generic ones.
- So actually, when you write

```
class Foo<T> {  
    T x;  
    T mogrify(T y) { ... }  
}  
  
Foo<Integer> q = new Foo<Integer>();  
Integer r = q.mogrify(s);
```

Java really gives you

```
class Foo {  
    Object x;  
    Object mogrify(Object y) { ... }  
}  
  
Foo q = new Foo();  
Integer r =  
(Integer) q.mogrify((Integer) s);
```

That is, it supplies the casts automatically, and also throws in some additional checks. If it can't guarantee that all those casts will work, gives you a warning about "unsafe" constructs.

Limitations

Because of Java's design choices, there are some limitations to generic programming:

- Since all kinds of Foo or List are really the same,
 - `L instanceof List<String>` will be true when `L` is a `List<Integer>`.
 - Inside, e.g., class Foo, you cannot write `new T()`, `new T[]`, or `x instanceof T`.
- Primitive types are not allowed as type parameters.
 - Can't have `ArrayList<int>`, just `ArrayList<Integer>`.
 - Fortunately, automatic boxing and unboxing makes this substitution easy:

```
int sum(ArrayList<Integer> L) {  
    int N;  N = 0;  
    for (int x : L) { N += x; }  
    return N;  
}
```

- Unfortunately, boxing and unboxing have significant costs.

CS61B Lecture #26

Today:

- Sorting algorithms: why?
- Insertion Sort.
- Inversions
- Selection sorting.
- Heapsort.

Purposes of Sorting

- Sorting supports searching; binary search is a standard example.
- Also supports other kinds of search:
 - Are there two equal items in this set?
 - Are there two items in this set that both have the same value for property X?
 - What are my nearest neighbors?
- Used in numerous unexpected algorithms, such as convex hull (smallest convex polygon enclosing set of points).

Some Definitions

- A *sorting algorithm* (or *sort*) *permutes* (re-arranges) a sequence of elements to brings them into order, according to some *total order*.
- A total order, \preceq , is:
 - **Total:** $x \preceq y$ or $y \preceq x$ for all x, y .
 - **Reflexive:** $x \preceq x$;
 - **Antisymmetric:** $x \preceq y$ and $y \preceq x$ iff $x = y$.
 - **Transitive:** $x \preceq y$ and $y \preceq z$ implies $x \preceq z$.
- However, our orderings may treat unequal items as equivalent:
 - E.g., there can be two dictionary definitions for the same word. If we sort only by the word being defined (ignoring the definition), then sorting could put either entry first.
 - A sort that does not change the relative order of equivalent entries (compared to the input) is called *stable*.

Classifications

- *Internal sorts* keep all data in primary memory.
- *External sorts* process large amounts of data in batches, keeping what won't fit in secondary storage (in the old days, tapes).
- *Comparison-based* sorting assumes only thing we know about keys is their order.
- *Radix sorting* uses more information about key structure.
- *Insertion sorting* works by repeatedly inserting items at their appropriate positions in the sorted sequence being constructed.
- *Selection sorting* works by repeatedly selecting the next larger (smaller) item in order and adding it to one end of the sorted sequence being constructed.

Sorting Arrays of Primitive Types in the Java Library

- The java library provides static methods to sort arrays in the class `java.util.Arrays`.
- For each primitive type P other than boolean, there are

```
/** Sort all elements of ARR into non-descending order. */
static void sort(P[] arr) { ... }

/** Sort elements FIRST .. END-1 of ARR into non-descending
 * order. */
static void sort(P[] arr, int first, int end) { ... }

/** Sort all elements of ARR into non-descending order,
 * possibly using multiprocessing for speed. */
static void parallelSort(P[] arr) { ... }

/** Sort elements FIRST .. END-1 of ARR into non-descending
 * order, possibly using multiprocessing for speed. */
static void parallelSort(P[] arr, int first, int end) {...}
```

Sorting Arrays of Reference Types in the Java Library

- For reference types, C, that have a *natural order* (that is, that implement `java.lang.Comparable`), we have four analogous methods (one-argument sort, three-argument sort, and two parallelSort methods):

```
static <C extends Comparable<? super C>> void sort(C[] arr) {...}  
static <C extends Comparable<? super C>>  
    void sort(C[] arr, int first, int end) { ... }  
static <C extends Comparable<? super C>>  
    void parallelSort(P[] arr) { ... }  
static <C extends Comparable<? super C>>  
    void parallelSort(P[] arr, int first, int end) {...}
```

- And for all reference types, R, we have four more:

```
/** Sort all elements of ARR stably into non-descending order  
 * according to the ordering defined by COMP. */  
static <R> void sort(R[] arr, Comparator<? super R> comp) {...}  
and so forth.
```

- Q: Why the fancy generic arguments for comp?

Sorting Arrays of Reference Types in the Java Library

- For reference types, C, that have a *natural order* (that is, that implement `java.lang.Comparable`), we have four analogous methods (one-argument sort, three-argument sort, and two parallelSort methods):

```
static <C extends Comparable<? super C>> void sort(C[] arr) {...}
static <C extends Comparable<? super C>>
    void sort(C[] arr, int first, int end) { ... }
static <C extends Comparable<? super C>>
    void parallelSort(P[] arr) { ... }
static <C extends Comparable<? super C>>
    void parallelSort(P[] arr, int first, int end) {...}
```

- And for all reference types, R, we have four more:

```
/** Sort all elements of ARR stably into non-descending order
 * according to the ordering defined by COMP. */
static <R> void sort(R[] arr, Comparator<? super R> comp) {...}
and so forth.
```

- Q: Why the fancy generic arguments for comp?

A: We want to allow types that have `compareTo` methods that apply also to more general types.

Sorting Lists in the Java Library

- The class `java.util.Collections` contains two methods similar to the sorting methods for arrays of reference types:

```
/** Sort all elements of LST stably into non-descending
 * order. */
static <C extends Comparable<? super C>> sort(List<C> lst) {...}
etc.
```

```
/** Sort all elements of LST stably into non-descending
 * order according to the ordering defined by COMP. */
static <R> void sort(List<R> , Comparator<? super R> comp) {...}
etc.
```

- Also a default instance method in the `List<R>` interface itself:

```
/** Sort all elements of LST stably into non-descending
 * order according to the ordering defined by COMP. */
default void sort(Comparator<? super R> comp) {...}
```

Examples

- Assume:

```
import static java.util.Arrays.*;  
import static java.util.Collections.*;
```

- Sort X, a String[] or List<String>, into non-descending order:

```
sort(X);    // or ...
```

- Sort X into reverse order (Java 8):

```
sort(X, (String x, String y) -> { return y.compareTo(x); });  
// or  
sort(X, Collections.reverseOrder());    // or  
X.sort(Collections.reverseOrder());      // for X a List
```

- Sort X[10], ..., X[100] in array or List X (rest unchanged):

```
sort(X, 10, 101);
```

- Sort L[10], ..., L[100] in list L (rest unchanged):

```
sort(L.sublist(10, 101));
```

Sorting by Insertion

- Simple idea:
 - Starting with an empty sequence of outputs.
 - Add each item from the input, *inserting* it into the output sequence at the right point.
- Very simple, good for small sets of data.
- With vector or linked list, time for find + insert of one item is at worst $\Theta(k)$, where k is # of outputs so far.
- This gives us a $\Theta(N^2)$ algorithm (worst case as usual).
- Can we say more?

Inversions

- Can run in $\Theta(N)$ comparisons if already sorted.
- Consider a typical implementation for arrays:

```
for (int i = 1; i < A.length; i += 1) {  
    int j;  
    Object x = A[i];  
    for (j = i-1; j >= 0; j -= 1) {  
        if (A[j].compareTo(x) <= 0) /* (1) */  
            break;  
        A[j+1] = A[j]; /* (2) */  
    }  
    A[j+1] = x;  
}
```

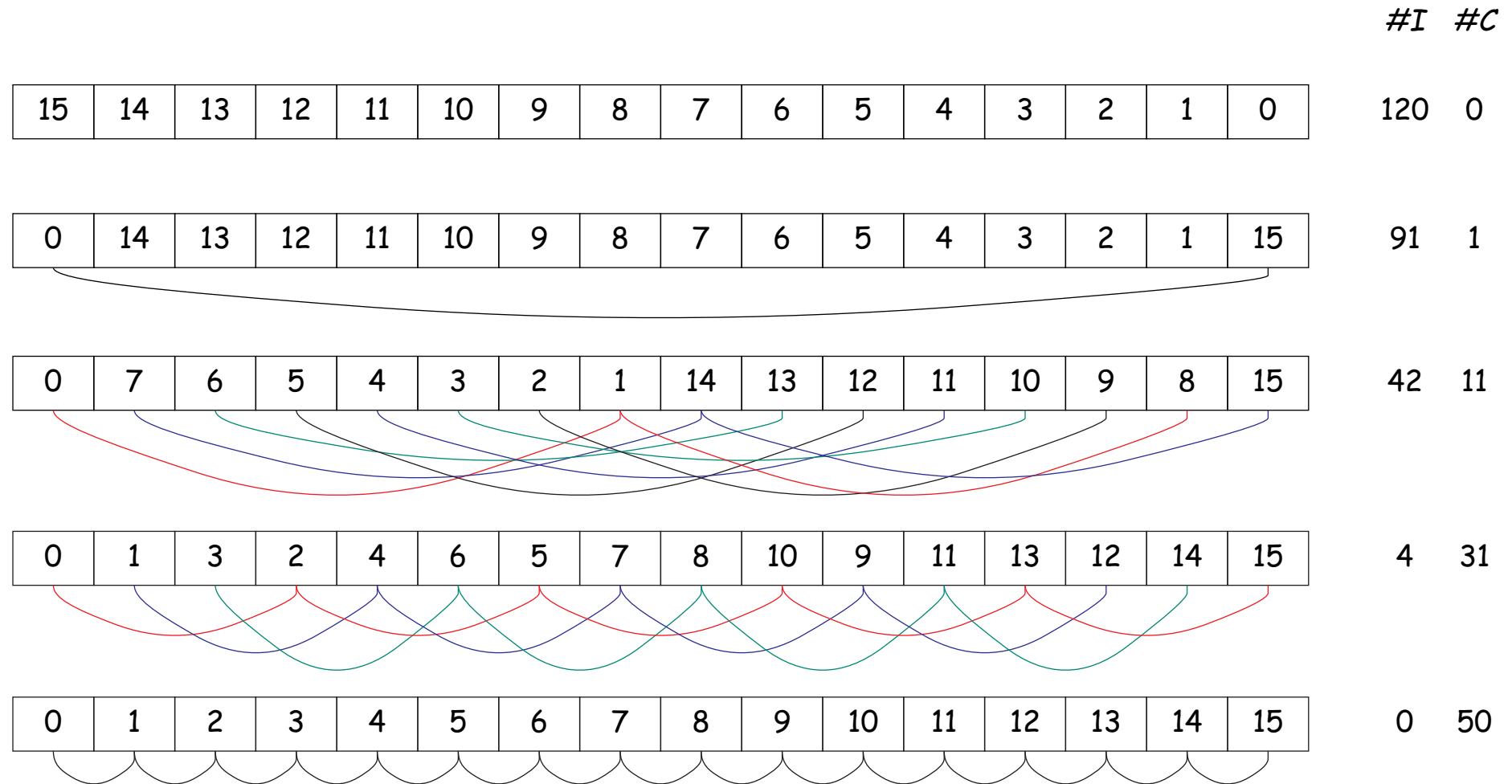
- #times (1) executes for each $j \approx$ how far x must move.
- If each item is within K of its proper places, requires $O(KN)$ operations.
- Thus is good for any amount of *nearly sorted* data.
- One measure of unsortedness: # of *inversions*: pairs that are out of order (= 0 when sorted, $N(N - 1)/2$ when reversed).
- Each execution of (2) decreases inversions by 1.

Shell's sort

Idea: Improve insertion sort by first sorting *distant* elements:

- First sort subsequences of elements $2^k - 1$ apart:
 - sort items #0, $2^k - 1, 2(2^k - 1), 3(2^k - 1), \dots$, then
 - sort items #1, $1 + 2^k - 1, 1 + 2(2^k - 1), 1 + 3(2^k - 1), \dots$, then
 - sort items #2, $2 + 2^k - 1, 2 + 2(2^k - 1), 2 + 3(2^k - 1), \dots$, then
 - etc.
 - sort items # $2^k - 2, 2(2^k - 1) - 1, 3(2^k - 1) - 1, \dots$,
 - Each time an item moves, can reduce #inversions by as much as $2^{k+1} - 3$.
- Now sort subsequences of elements $2^{k-1} - 1$ apart:
 - sort items #0, $2^{k-1} - 1, 2(2^{k-1} - 1), 3(2^{k-1} - 1), \dots$, then
 - sort items #1, $1 + 2^{k-1} - 1, 1 + 2(2^{k-1} - 1), 1 + 3(2^{k-1} - 1), \dots$,
 - :
- End at plain insertion sort ($2^0 = 1$ apart), but with most inversions gone.
- Sort is $\Theta(N^{3/2})$ (take CS170 for why!).

Example of Shell's Sort



I: Inversions left.

C: Cumulative comparisons used to sort subsequences by insertion sort.

Sorting by Selection: Straight Selection

- Idea: Keep selecting the smallest (or largest) element, and adding it to the appropriate end of the result.
- Gives us *straight selection sort*:

```
<T extends Comparable<T>> void selectionSort(T[] A) {  
    for (int i = A.length; i > 0; i --) {  
        int maxIndex;  
        maxIndex = 0;  
        for (int k = 1; k < i; k += 1) {  
            if (A[maxIndex].compareTo(A[k]) < 0) {  
                maxIndex = k;  
            }  
        }  
        T temp = A[maxIndex]; A[maxIndex] = A[i - 1]; A[i - 1] = temp;  
    }  
}
```

- Time bound?

Sorting by Selection: Straight Selection

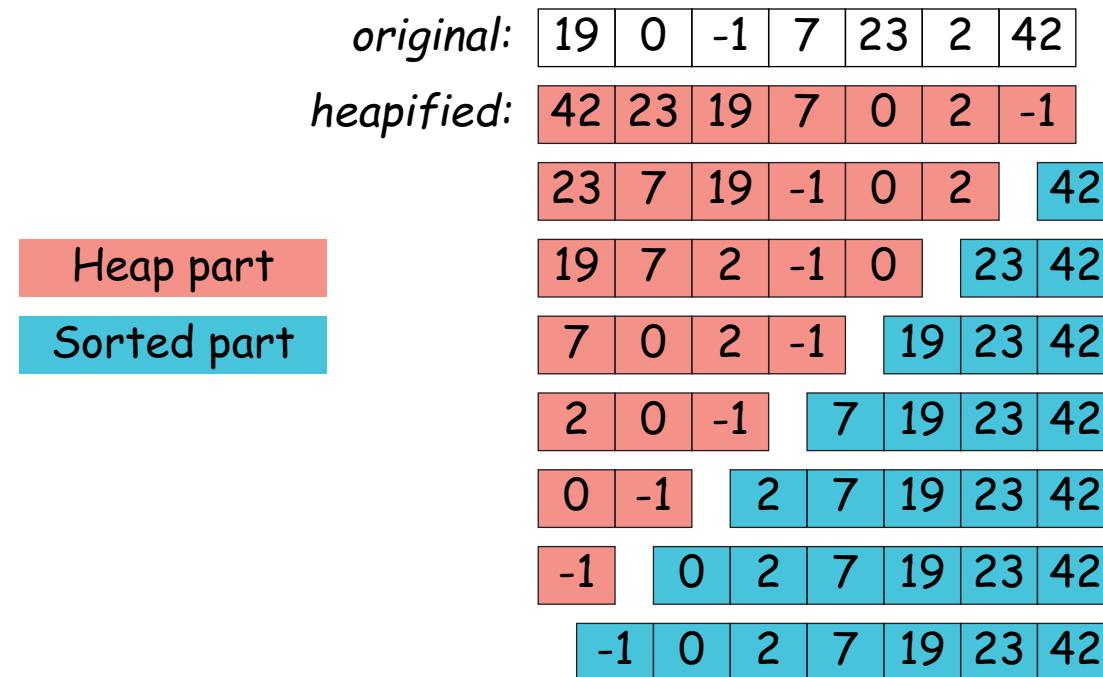
- Idea: Keep selecting the smallest (or largest) element, and adding it to the appropriate end of the result.
- Gives us *straight selection sort*:

```
<T extends Comparable<T>> void selectionSort(T[] A) {  
    for (int i = A.length; i > 0; i --) {  
        int maxIndex;  
        maxIndex = 0;  
        for (int k = 1; k < i; k += 1) {  
            if (A[maxIndex].compareTo(A[k]) < 0) {  
                maxIndex = k;  
            }  
        }  
        T temp = A[maxIndex]; A[maxIndex] = A[i - 1]; A[i - 1] = temp;  
    }  
}
```

- Clearly this is $\Theta(N^2)$ in the worst case.
- But it is also $\Theta(N^2)$ in the best case, regardless of how sorted the data originally are.

Sorting by Selection: Heapsort

- So straight selection is not such a great idea on a simple list or vector.
- But we've already seen selection in action elsewhere: use heap.
- Gives $O(N \lg N)$ algorithm (N remove-first operations).
- And since we remove items from the end of the heap, we can use that area to accumulate the result:



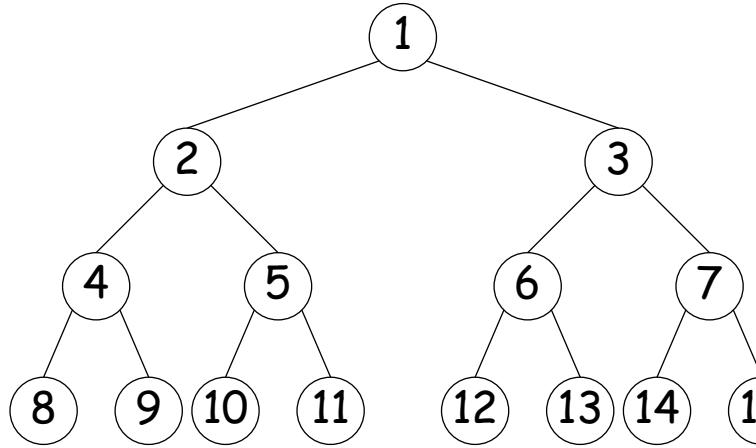
Sorting By Selection: Initial Heapifying

- When covering heaps before, we created them by insertion in an initially empty heap.
- When given an array of unheaped data to start with, there is a faster procedure (assume heap indexed from 0):

```
void heapify(int[] arr) {  
    int N = arr.length;  
    for (int k = N / 2; k >= 0; k -= 1) {  
        for (int p = k, c = 0; 2*p + 1 < N; p = c) {  
            reheapify downward from p;  
        }  
    }  
}
```

- At each iteration of the p loop, only the element at p might be out of order with respect to its descendants, so reheapifying downward will restore the subtree rooted at p to proper heap ordering.
- Looks like the procedure for re-inserting an element after the top element of the heap is removed, repeated $N/2$ times, *but,...*

Cost of Creating Heap



1 node \times 3 steps down

2 nodes \times 2 steps down

4 nodes \times 1 step down

- In general, worst-case cost for a heap with $h + 1$ levels is

$$\begin{aligned} & 2^0 \cdot h + 2^1 \cdot (h - 1) + \dots + 2^{h-1} \cdot 1 \\ &= (2^0 + 2^1 + \dots + 2^{h-1}) + (2^0 + 2^1 + \dots + 2^{h-2}) + \dots + (2^0) \\ &= (2^h - 1) + (2^{h-1} - 1) + \dots + (2^1 - 1) \\ &= 2^{h+1} - 1 - h \\ &\in \Theta(2^h) = \Theta(N) : \text{linear time, not } N \lg N. \end{aligned}$$

- Alas, since the rest of heapsort still takes $\Theta(N \lg N)$, this does not improve its overall asymptotic cost.

CS61B Lectures #27

Today:

- Merge sorts
- Quicksort

Readings: Today: *DS(IJ)*, Chapter 8; Next topic: Chapter 9.

Merge Sorting

Idea: Divide data into subsequences; recursively sort the subsequences; merge results.

- We've already seen the analysis (Lecture #16): $\Theta(N \lg N)$.
- Good for *external sorting*:
 - First break the data into small enough chunks to fit in memory and sort each.
 - Then repeatedly merge into bigger and bigger sequences.
- Can merge K sorted sequences of *arbitrary size* on secondary storage using $\Theta(K)$ storage:

Data[] V = new Data[K] ;

For all i, set V[i] to the first data item of sequence i;
while there is data left to sort:

 Find k so that V[k] has data and is smallest;

 Write V[k] to the output sequence;

 If there is more data in sequence k, read it into V[k],
 otherwise, clear V[k] ;

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

55	20	31	80	58
----	----	----	----	----

Input subsequences:

25	-4	34	16	8
----	----	----	----	---

61	39	35	42	60
----	----	----	----	----

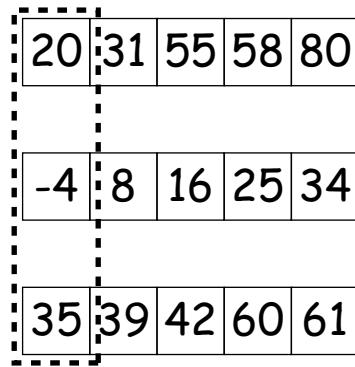
- First, the input data sequence is divided into subsequences.

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Sorted subsequences:



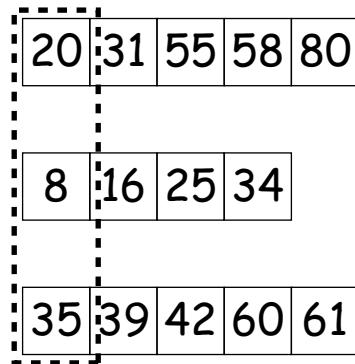
- First, the input data sequence is divided into subsequences.
- Next, the subsequences are themselves sorted (possibly by a recursive application of mergesort.)

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Remaining subsequences:



Result:

-4

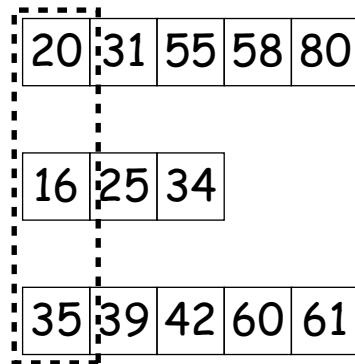
- The dashed window shows the input data that must be in memory at any given time. It is of constant size, no matter what the size of the input.
- One by one, the smallest item in the dashed window is removed from its sequence and added to the result.

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Remaining subsequences:



Result:

-4	8
----	---

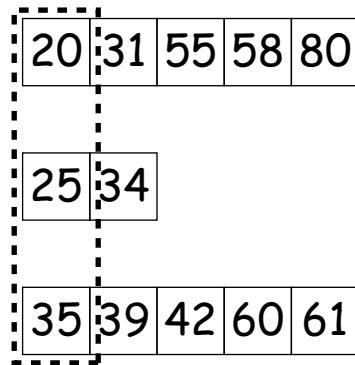
- The dashed window shows the input data that must be in memory at any given time. It is of constant size, no matter what the size of the input.
- One by one, the smallest item in the dashed window is removed from its sequence and added to the result.

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Remaining subsequences:



Result:

-4	8	16
----	---	----

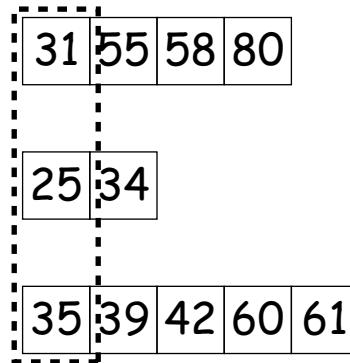
- The dashed window shows the input data that must be in memory at any given time. It is of constant size, no matter what the size of the input.
- One by one, the smallest item in the dashed window is removed from its sequence and added to the result.

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Remaining subsequences:



Result:

-4	8	16	20
----	---	----	----

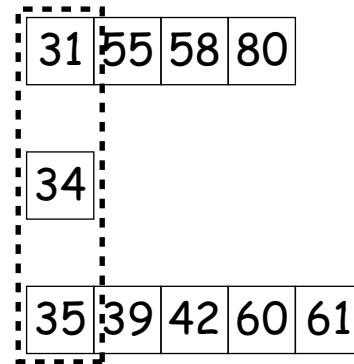
- The dashed window shows the input data that must be in memory at any given time. It is of constant size, no matter what the size of the input.
- One by one, the smallest item in the dashed window is removed from its sequence and added to the result.

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Remaining subsequences:



Result:

-4	8	16	20	25
----	---	----	----	----

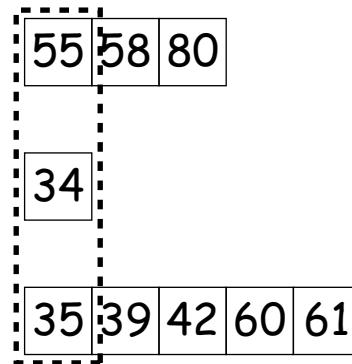
- The dashed window shows the input data that must be in memory at any given time. It is of constant size, no matter what the size of the input.
- One by one, the smallest item in the dashed window is removed from its sequence and added to the result.

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Remaining subsequences:



Result:

-4	8	16	20	25	31
----	---	----	----	----	----

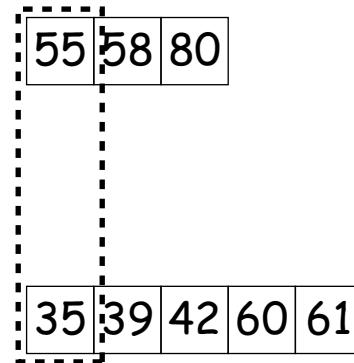
- The dashed window shows the input data that must be in memory at any given time. It is of constant size, no matter what the size of the input.
- One by one, the smallest item in the dashed window is removed from its sequence and added to the result.

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Remaining subsequences:



Result:

-4	8	16	20	25	31	34
----	---	----	----	----	----	----

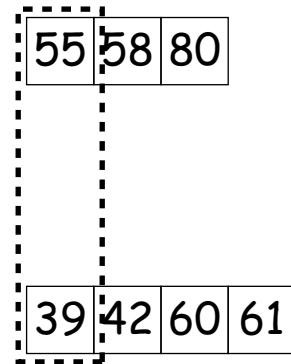
- The dashed window shows the input data that must be in memory at any given time. It is of constant size, no matter what the size of the input.
- One by one, the smallest item in the dashed window is removed from its sequence and added to the result.

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Remaining subsequences:



Result:

-4	8	16	20	25	31	34	35
----	---	----	----	----	----	----	----

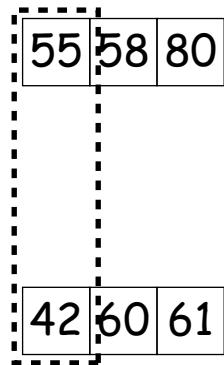
- The dashed window shows the input data that must be in memory at any given time. It is of constant size, no matter what the size of the input.
- One by one, the smallest item in the dashed window is removed from its sequence and added to the result.

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Remaining subsequences:



Result:

-4	8	16	20	25	31	34	35	39
----	---	----	----	----	----	----	----	----

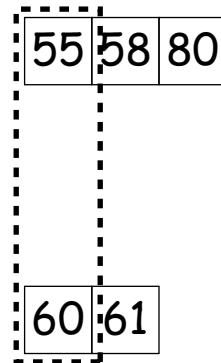
- The dashed window shows the input data that must be in memory at any given time. It is of constant size, no matter what the size of the input.
- One by one, the smallest item in the dashed window is removed from its sequence and added to the result.

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Remaining subsequences:



Result:

-4	8	16	20	25	31	34	35	39	42
----	---	----	----	----	----	----	----	----	----

- The dashed window shows the input data that must be in memory at any given time. It is of constant size, no matter what the size of the input.
- One by one, the smallest item in the dashed window is removed from its sequence and added to the result.

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Remaining subsequences:



Result:

-4	8	16	20	25	31	34	35	39	42	55
----	---	----	----	----	----	----	----	----	----	----

- The dashed window shows the input data that must be in memory at any given time. It is of constant size, no matter what the size of the input.
- One by one, the smallest item in the dashed window is removed from its sequence and added to the result.

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Remaining subsequences:



Result:

-4	8	16	20	25	31	34	35	39	42	55	58
----	---	----	----	----	----	----	----	----	----	----	----

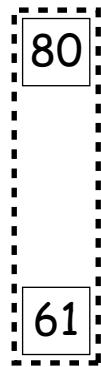
- The dashed window shows the input data that must be in memory at any given time. It is of constant size, no matter what the size of the input.
- One by one, the smallest item in the dashed window is removed from its sequence and added to the result.

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Remaining subsequences:



Result:

-4	8	16	20	25	31	34	35	39	42	55	58	60
----	---	----	----	----	----	----	----	----	----	----	----	----

- The dashed window shows the input data that must be in memory at any given time. It is of constant size, no matter what the size of the input.
- One by one, the smallest item in the dashed window is removed from its sequence and added to the result.

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Remaining subsequences:



Result:

-4	8	16	20	25	31	34	35	39	42	55	58	60	61
----	---	----	----	----	----	----	----	----	----	----	----	----	----

- The dashed window shows the input data that must be in memory at any given time. It is of constant size, no matter what the size of the input.
- One by one, the smallest item in the dashed window is removed from its sequence and added to the result.

Merge Sorting Illustrated

Input:

55	20	31	80	58	25	-4	34	16	8	61	39	35	42	60
----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

Remaining subsequences:



Result:

-4	8	16	20	25	31	34	35	39	42	55	58	60	61	80
----	---	----	----	----	----	----	----	----	----	----	----	----	----	----

- The dashed window shows the input data that must be in memory at any given time. It is of constant size, no matter what the size of the input.
- One by one, the smallest item in the dashed window is removed from its sequence and added to the result.

Illustration of Internal Merge Sort

For internal sorting, we can use a *binomial comb* to orchestrate an iterative merge sort.

- Start with $\lg N + 1$ buckets that can contain sublists, initially empty.
- Bucket # k is either empty or contains 2^k sorted items at any time.
- For each item in the input list, turn it into a 1-element list, and merge it into bucket 0 (or simply put it in bucket 0 if that is empty).
- You will only merge lists of length 2^k into bucket k . Whenever that gives a list of size 2^{k+1} , merge it into bucket $k + 1$ and clear bucket k (and so on as needed with buckets $k + 2$, etc.)
- When all inputs are processed, merge all the buckets into the final list.

L: (9, 15, 5, 3, 0, 6, 10, -1, 2, 20, 8)

0:	0
1:	0
2:	0
3:	0

Merge (9)

Illustration of Internal Merge Sort

For internal sorting, we can use a *binomial comb* to orchestrate an iterative merge sort.

- Start with $\lg N + 1$ buckets that can contain sublists, initially empty.
- Bucket # k is either empty or contains 2^k sorted items at any time.
- For each item in the input list, turn it into a 1-element list, and merge it into bucket 0 (or simply put it in bucket 0 if that is empty).
- You will only merge lists of length 2^k into bucket k . Whenever that gives a list of size 2^{k+1} , merge it into bucket $k + 1$ and clear bucket k (and so on as needed with buckets $k + 2$, etc.)
- When all inputs are processed, merge all the buckets into the final list.

L: (9, 15, 5, 3, 0, 6, 10, -1, 2, 20, 8)

0:	1	●	→	(9)
1:	0			
2:	0			
3:	0			

Illustration of Internal Merge Sort

For internal sorting, we can use a *binomial comb* to orchestrate an iterative merge sort.

- Start with $\lg N + 1$ buckets that can contain sublists, initially empty.
- Bucket # k is either empty or contains 2^k sorted items at any time.
- For each item in the input list, turn it into a 1-element list, and merge it into bucket 0 (or simply put it in bucket 0 if that is empty).
- You will only merge lists of length 2^k into bucket k . Whenever that gives a list of size 2^{k+1} , merge it into bucket $k + 1$ and clear bucket k (and so on as needed with buckets $k + 2$, etc.)
- When all inputs are processed, merge all the buckets into the final list.

L: (9, 15, 5, 3, 0, 6, 10, -1, 2, 20, 8)

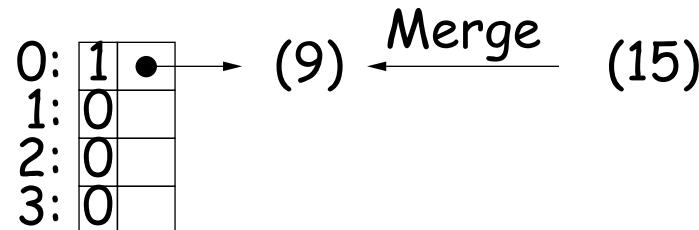


Illustration of Internal Merge Sort

For internal sorting, we can use a *binomial comb* to orchestrate an iterative merge sort.

- Start with $\lg N + 1$ buckets that can contain sublists, initially empty.
- Bucket # k is either empty or contains 2^k sorted items at any time.
- For each item in the input list, turn it into a 1-element list, and merge it into bucket 0 (or simply put it in bucket 0 if that is empty).
- You will only merge lists of length 2^k into bucket k . Whenever that gives a list of size 2^{k+1} , merge it into bucket $k + 1$ and clear bucket k (and so on as needed with buckets $k + 2$, etc.)
- When all inputs are processed, merge all the buckets into the final list.

L: (9, 15, 5, 3, 0, 6, 10, -1, 2, 20, 8)

0:	0
1:	0
2:	0
3:	0

$\xleftarrow{\text{Merge}}$ (9, 15)

Illustration of Internal Merge Sort

For internal sorting, we can use a *binomial comb* to orchestrate an iterative merge sort.

- Start with $\lg N + 1$ buckets that can contain sublists, initially empty.
- Bucket # k is either empty or contains 2^k sorted items at any time.
- For each item in the input list, turn it into a 1-element list, and merge it into bucket 0 (or simply put it in bucket 0 if that is empty).
- You will only merge lists of length 2^k into bucket k . Whenever that gives a list of size 2^{k+1} , merge it into bucket $k + 1$ and clear bucket k (and so on as needed with buckets $k + 2$, etc.)
- When all inputs are processed, merge all the buckets into the final list.

L: (9, 15, 5, 3, 0, 6, 10, -1, 2, 20, 8)

0:	0
1:	1 •
2:	0
3:	0

→ (9, 15)

Illustration of Internal Merge Sort

For internal sorting, we can use a *binomial comb* to orchestrate an iterative merge sort.

- Start with $\lg N + 1$ buckets that can contain sublists, initially empty.
- Bucket # k is either empty or contains 2^k sorted items at any time.
- For each item in the input list, turn it into a 1-element list, and merge it into bucket 0 (or simply put it in bucket 0 if that is empty).
- You will only merge lists of length 2^k into bucket k . Whenever that gives a list of size 2^{k+1} , merge it into bucket $k + 1$ and clear bucket k (and so on as needed with buckets $k + 2$, etc.)
- When all inputs are processed, merge all the buckets into the final list.

L: (9, 15, 5, 3, 0, 6, 10, -1, 2, 20, 8)

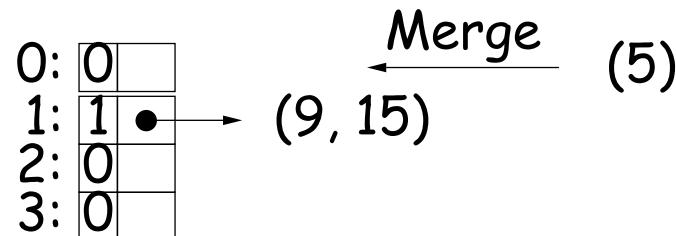


Illustration of Internal Merge Sort

For internal sorting, we can use a *binomial comb* to orchestrate an iterative merge sort.

- Start with $\lg N + 1$ buckets that can contain sublists, initially empty.
- Bucket # k is either empty or contains 2^k sorted items at any time.
- For each item in the input list, turn it into a 1-element list, and merge it into bucket 0 (or simply put it in bucket 0 if that is empty).
- You will only merge lists of length 2^k into bucket k . Whenever that gives a list of size 2^{k+1} , merge it into bucket $k + 1$ and clear bucket k (and so on as needed with buckets $k + 2$, etc.)
- When all inputs are processed, merge all the buckets into the final list.

L: (9, 15, 5, 3, 0, 6, 10, -1, 2, 20, 8)

0:	<table border="1"><tr><td>1</td><td>•</td></tr></table>	1	•	→ (5)
1	•			
1:	<table border="1"><tr><td>1</td><td>•</td></tr></table>	1	•	→ (9, 15)
1	•			
2:	<table border="1"><tr><td>0</td></tr></table>	0		
0				
3:	<table border="1"><tr><td>0</td></tr></table>	0		
0				

Illustration of Internal Merge Sort

For internal sorting, we can use a *binomial comb* to orchestrate an iterative merge sort.

- Start with $\lg N + 1$ buckets that can contain sublists, initially empty.
- Bucket # k is either empty or contains 2^k sorted items at any time.
- For each item in the input list, turn it into a 1-element list, and merge it into bucket 0 (or simply put it in bucket 0 if that is empty).
- You will only merge lists of length 2^k into bucket k . Whenever that gives a list of size 2^{k+1} , merge it into bucket $k + 1$ and clear bucket k (and so on as needed with buckets $k + 2$, etc.)
- When all inputs are processed, merge all the buckets into the final list.

L: (9, 15, 5, 3, 0, 6, 10, -1, 2, 20, 8)

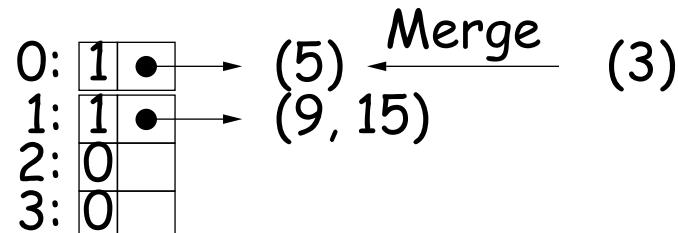


Illustration of Internal Merge Sort

For internal sorting, we can use a *binomial comb* to orchestrate an iterative merge sort.

- Start with $\lg N + 1$ buckets that can contain sublists, initially empty.
- Bucket # k is either empty or contains 2^k sorted items at any time.
- For each item in the input list, turn it into a 1-element list, and merge it into bucket 0 (or simply put it in bucket 0 if that is empty).
- You will only merge lists of length 2^k into bucket k . Whenever that gives a list of size 2^{k+1} , merge it into bucket $k + 1$ and clear bucket k (and so on as needed with buckets $k + 2$, etc.)
- When all inputs are processed, merge all the buckets into the final list.

L: (9, 15, 5, 3, 0, 6, 10, -1, 2, 20, 8)

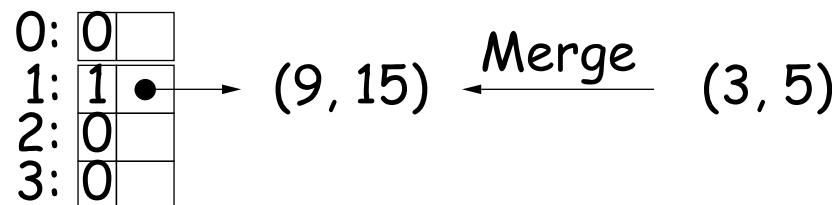


Illustration of Internal Merge Sort

For internal sorting, we can use a *binomial comb* to orchestrate an iterative merge sort.

- Start with $\lg N + 1$ buckets that can contain sublists, initially empty.
- Bucket # k is either empty or contains 2^k sorted items at any time.
- For each item in the input list, turn it into a 1-element list, and merge it into bucket 0 (or simply put it in bucket 0 if that is empty).
- You will only merge lists of length 2^k into bucket k . Whenever that gives a list of size 2^{k+1} , merge it into bucket $k + 1$ and clear bucket k (and so on as needed with buckets $k + 2$, etc.)
- When all inputs are processed, merge all the buckets into the final list.

L: (9, 15, 5, 3, 0, 6, 10, -1, 2, 20, 8)

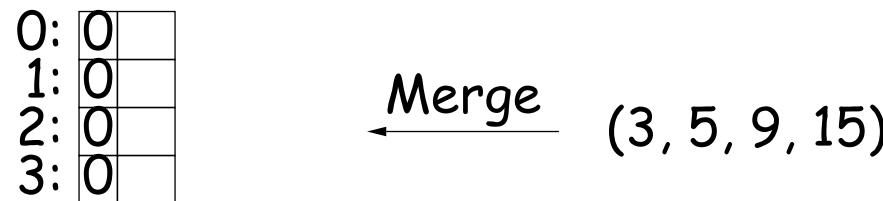


Illustration of Internal Merge Sort

For internal sorting, we can use a *binomial comb* to orchestrate an iterative merge sort.

- Start with $\lg N + 1$ buckets that can contain sublists, initially empty.
- Bucket # k is either empty or contains 2^k sorted items at any time.
- For each item in the input list, turn it into a 1-element list, and merge it into bucket 0 (or simply put it in bucket 0 if that is empty).
- You will only merge lists of length 2^k into bucket k . Whenever that gives a list of size 2^{k+1} , merge it into bucket $k + 1$ and clear bucket k (and so on as needed with buckets $k + 2$, etc.)
- When all inputs are processed, merge all the buckets into the final list.

L: (9, 15, 5, 3, 0, 6, 10, -1, 2, 20, 8)

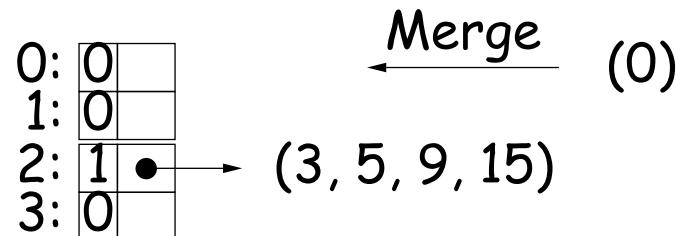


Illustration of Internal Merge Sort (II)

L: (9, 15, 5, 3, 0, 6, 10, -1, 2, 20, 8)

0:	0	
1:	0	
2:	0	
3:	0	

0 elements processed

0:	1	•	→ (9)
1:	0		
2:	0		
3:	0		

1 element processed

0:	0		→ (9, 15)
1:	1	•	
2:	0		
3:	0		

2 elements processed

0:	1	•	→ (5)
1:	1	•	→ (9, 15)
2:	0		
3:	0		

3 elements processed

0:	0		
1:	0		
2:	1	•	→ (3, 5, 9, 15)
3:	0		

4 elements processed

0:	0		→ (0, 6)
1:	1	•	
2:	1	•	→ (3, 5, 9, 15)
3:	0		

6 elements processed

0:	1	•	→ (8)
1:	1	•	→ (2, 20)
2:	0		
3:	1	•	→ (-1, 0, 3, 5, 6, 9, 10, 15)

11 elements processed

Final Step: Merge all the lists into (-1, 0, 2, 3, 5, 6, 8, 9, 10, 15, 20)

Quicksort: Speed through Probability

Idea:

- *Partition* data into pieces: everything $>$ a *pivot* value at the high end of the sequence to be sorted, and everything $<$ on the low end, and everything $=$ between.
- Repeat recursively on the high and low pieces.
- For speed, stop when pieces are “small enough” and do *insertion sort* on the whole thing.
- Reason: *insertion sort* has low constant factors. By design, no item will move out of its piece [why?], so when pieces are small, #inversions is, too.
- Have to choose pivot well. E.g.: *median* of first, last and middle items of sequence.

Example of Quicksort

- In this example, we continue until pieces are size ≤ 4 .
- Pivots for next step are starred. We arrange to move the pivot to dividing line each time.
- Last step is insertion sort.

16	10	13	18	-4	-7	12	-5	19	15	0	22	29	34	-1*
-4	-5	-7	-1	18	13	12	10	19	15	0	22	29	34	16*
-4	-5	-7	-1	15	13	12*	10	0	16	19*	22	29	34	18
-4	-5	-7	-1	10	0	12	15	13	16	18	19	29	34	22

- Now everything is “close to” right (just 7 inversions), so just do insertion sort:

-7	-5	-4	-1	0	10	12	13	15	16	18	19	22	29	34
----	----	----	----	---	----	----	----	----	----	----	----	----	----	----

Performance of Quicksort

- Probabalistic time:
 - With a good choice of pivots, we divide data by two each time, giving $\Theta(N \lg N)$ and a good constant factor relative to merge or heap sort.
 - With a bad choice of pivots, most items will be on one side each time: leading to a $\Theta(N^2)$ time.
 - Time is $\Omega(N \lg N)$ even in the best case, so insertion sort is better for nearly ordered input sets.
- Interesting point: randomly shuffling the data before sorting makes $\Omega(N^2)$ time **very** unlikely!

Quick Selection

The Selection Problem: for given k , find k^{th} smallest element in data.

- Obvious method: sort, select element # k , time $\Theta(N \lg N)$.
- If $k \leq$ some constant, we can easily do in $\Theta(N)$ time:
 - Go through array, keeping the smallest k items.
- Get *probably $\Theta(N)$ time* for all k by adapting quicksort:
 - Partition around some pivot, p , as in quicksort, arrange for that pivot to end up at the dividing line.
 - Suppose that in the result, the pivot is at index m , all elements \leq pivot have indices $\leq m$.
 - If $m = k$, you're done: p is answer.
 - If $m > k$, recursively select the k^{th} largest from the left half of the sequence.
 - If $m < k$, recursively select the $(k - m - 1)^{\text{th}}$ largest from the right half of sequence.

Selection Example

Problem: Find just item #10 in the sorted version of array:

Initial contents:

51	60	21	-4	37	4	49	10	40*	59	0	13	2	39	11	46	31
0																

Looking for #10 to left of pivot 40:

13	31	21	-4	37	4*	11	10	39	2	0	40	59	51	49	46	60
0																

Looking for #6 to right of pivot 4:

-4	0	2	4	37	13	11	10	39	21	31*	40	59	51	49	46	60
4																

Looking for #1 to right of pivot 31:

-4	0	2	4	21	13	11	10	31	39	37	40	59	51	49	46	60
9																

Just two elements; just sort and return #1:

-4	0	2	4	21	13	11	10	31	37	39	40	59	51	49	46	60
9																

Result: 39

Selection Performance

- For this algorithm, if m is roughly in the middle each time, the cost is

$$\begin{aligned} C(N) &= \begin{cases} 1, & \text{if } N = 1, \\ N + C(N/2), & \text{otherwise.} \end{cases} \\ &= N + N/2 + \dots + 1 \\ &= 2N - 1 \in \Theta(N) \end{aligned}$$

- But in worst case, we get $\Theta(N^2)$, as for quicksort.
- By another, non-obvious algorithm, we can get $\Theta(N)$ worst-case time for all k (take CS170).

CS61B Lectures #28

Today:

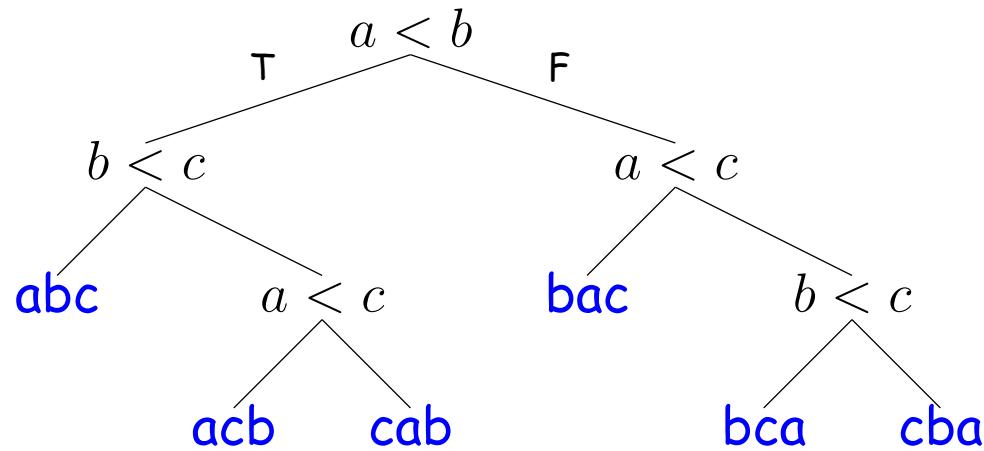
- Lower bounds on sorting by comparison
- Distribution counting, radix sorts

Readings: Today: *DS(IJ)*, Chapter 8; Next topic: Chapter 9.

Better than $N \lg N$?

- We can prove that *if all you can do to keys is compare them*, then sorting must take $\Omega(N \lg N)$.
- Basic idea: there are $N!$ possible ways the input data could be scrambled.
- Therefore, your program must be prepared to do $N!$ different combinations of data-moving operations.
- Therefore, there must be $N!$ possible combinations of outcomes of all the if-tests in your program, since those determine what data gets moved where (we're assuming that comparisons are 2-way).

Decision Tree
Height \propto Sorting time



Necessary Choices

- Since each if-test goes two ways, the number of possible different outcomes for k if-tests is 2^k .
- Thus, there must be enough tests so that $2^k \geq N!$, which means $k \geq \lg N!$.
- Using Stirling's approximation,

$$\begin{aligned}N! &\in \sqrt{2\pi N} \left(\frac{N}{e} \right)^N \left(1 + \Theta \left(\frac{1}{N} \right) \right), \\ \lg(N!) &\in \frac{1}{2}(\lg 2\pi + \lg N) + N \lg N - N \lg e + \lg \left(1 + \Theta \left(\frac{1}{N} \right) \right) \\ &= \Theta(N \lg N)\end{aligned}$$

- This tells us that k , the worst-case number of tests needed to sort N items by comparison sorting, is in $\Omega(N \lg N)$: there must be cases where we need (some multiple of) $N \lg N$ comparisons to sort N things.

Beyond Comparison: Distribution

- But suppose we can do more than compare keys?
- For example, how can we sort a set of N different integer keys whose values range from 0 to kN , for some small constant k ?
- One technique is *distribution sorting*:
 - Put the integers into N buckets; integer p goes to bucket $\lfloor p/k \rfloor$.
 - There are at most k keys per bucket, so concatenate and use insertion sort, which will now be fast.
- E.g., $k = 2, N = 10$:

Start:

14 3 10 13 4 2 19 17 0 9

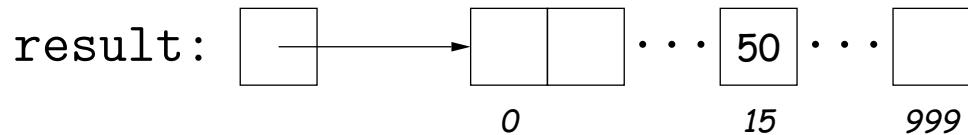
In buckets:

| 0 | 3 2 | 4 | | 9 | 10 | 13 | 14 | 17 | 19 |

- Now insertion sort is fast. Putting the data in buckets takes time $\Theta(N)$, and insertion sort takes $\Theta(kN)$. When k is fixed (constant), we have sorting in time $\Theta(N)$.

Distribution Counting

- Another technique: **count** the number of items $< 1, < 2$, etc.
- If $M_p = \#\text{items with value } < p$, then in sorted order, the j^{th} item with value p must be item $\#M_p + j$.
- Suppose that one has a set of numbers in the range $[0, 1000]$ (possibly with duplicates) and that exactly 15 of them are less than 50, which is also in the set. Then the result of sorting will look like this:

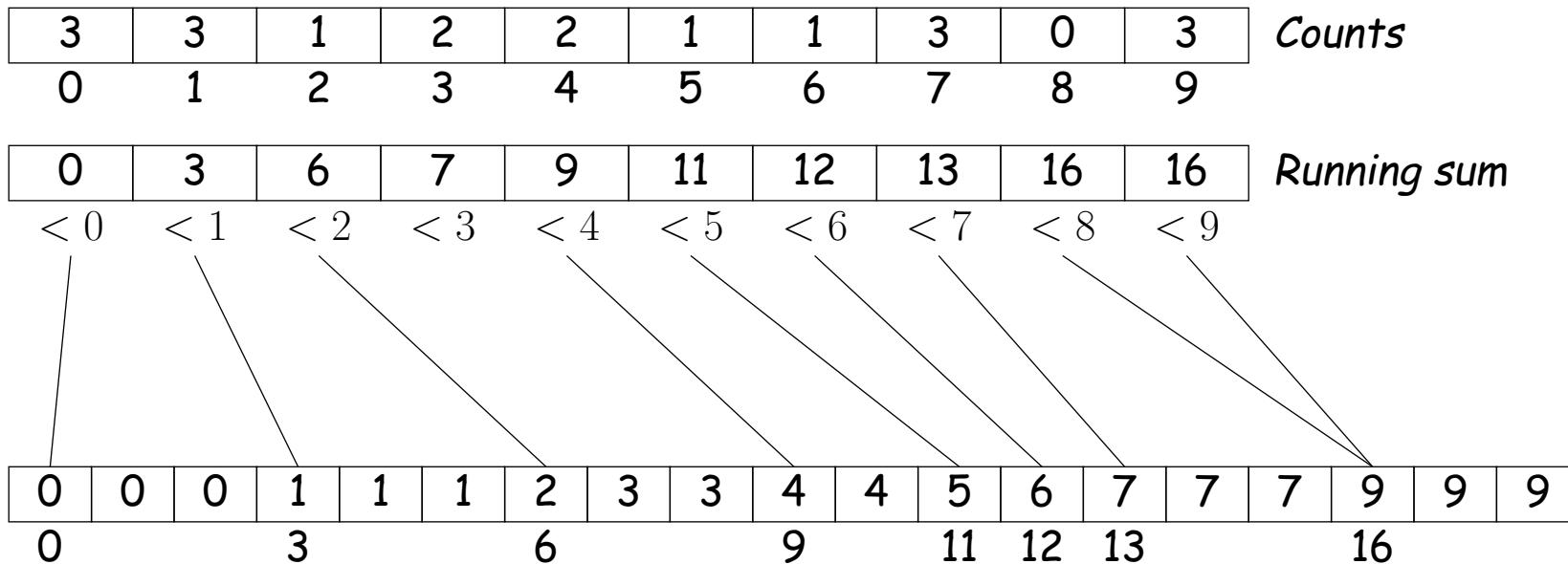


- In other words, the count of numbers $< k$ gives the index of k in the output array.
- If there are N items in the range $0..M - 1$, this gives another **linear-time**— $\Theta(M + N)$ —algorithm (We include M and N here to allow for both duplicates and for cases where $M \gg N$.)
- [Postscript on notation: the notations $[A, B]$, (A, B) , $[A, B)$, and $(A, B]$ above refer to **intervals**. The use of parentheses vs. square brackets reflects the distinction between **open** and **closed** intervals. Thus $x \in [A, B]$ iff $A \leq x \leq B$, while $x \in [A, B)$ iff $A \leq x < B$, etc.]

Distribution Counting Example

- Suppose all items are between 0 and 9 as in this example:

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



- "Counts" line gives # occurrences of each key.
- "Running sum" gives cumulative count of keys $<$ each value...
- ... which tells us where to put each key:
- The first instance of key k goes into slot m , where m is the number of key instances that are $< k$; next k goes into slot $m + 1$, etc.

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

0	3	6	7	9	11	12	13	16	16	Next positions
0	1	2	3	4	5	6	7	8	9	

																			Output
0																			

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

0	3	6	7	9	11	12	14	16	16	Next positions
0	1	2	3	4	5	6	7	8	9	

																			Output
0																			

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

1	3	6	7	9	11	12	14	16	16	Next positions
0	1	2	3	4	5	6	7	8	9	

0																			Output
0																			

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

1	3	6	7	10	11	12	14	16	16	Next positions
0	1	2	3	4	5	6	7	8	9	

0								4										Output
0								9										

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

2	3	6	7	10	11	12	14	16	16	Next positions
0	1	2	3	4	5	6	7	8	9	

0	0							4					7					Output
0		3		6				9					12				15	18

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

2	3	6	7	10	11	12	14	16	17	Next positions
0	1	2	3	4	5	6	7	8	9	

0	0						4					7			9			Output
0		3		6			9					12			15		18	

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

2	4	6	7	10	11	12	14	16	17	Next positions
0	1	2	3	4	5	6	7	8	9	

0	0		1			6			4				7			15	9		Output
0			3			6			9				12			18			

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

2	4	6	7	10	11	12	14	16	18	Next positions
0	1	2	3	4	5	6	7	8	9	

0	0		1			6			4				7			12		15	9	9		Output
0			3			6			9											18		

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

2	5	6	7	10	11	12	14	16	18	Next positions
0	1	2	3	4	5	6	7	8	9	

0	0		1	1				4				7			15	9	9		Output
0			3		6			9				12				15	9	9	

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

2	5	6	7	10	11	12	14	16	19	Next positions
0	1	2	3	4	5	6	7	8	9	

0	0		1	1				4				7			15	9	9	9	Output
0			3		6			9				12				15	18		

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

2	5	6	7	10	12	12	14	16	19	Next positions
0	1	2	3	4	5	6	7	8	9	

0	0		1	1			4		5	7		15	9	9	9	Output
0			3		6		9		12			18				

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

2	5	6	8	10	12	12	14	16	19	Next positions
0	1	2	3	4	5	6	7	8	9	

0	0		1	1		6	3		4		5		7			9	9	9	Output
0			3			6			9				12			15	18		

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

2	5	6	8	10	12	12	15	16	19	Next positions
0	1	2	3	4	5	6	7	8	9	

0	0		1	1		6	3		4		5		7	7		9	9	9	Output
0			3			6			9		12		15		18				

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

2	5	6	9	10	12	12	15	16	19	Next positions
0	1	2	3	4	5	6	7	8	9	

0	0		1	1		6	3	3	4		5		7	7	15	9	9	9	Output
0			3			6			9				12			15	18		

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

2	6	6	9	10	12	12	15	16	19	Next positions
0	1	2	3	4	5	6	7	8	9	

0	0		1	1	1	6	3	3	4		5	12	7	7	15	9	9	9	Output
0			3			6			9						15	18			

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

2	6	6	9	10	12	13	15	16	19	Next positions
0	1	2	3	4	5	6	7	8	9	

0	0		1	1	1	6	3	3	4		5	6	7	7	15	9	9	9	Output
0			3			6			9			12			15		18		

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

2	6	6	9	10	12	13	16	16	19	Next positions
0	1	2	3	4	5	6	7	8	9	

0	0		1	1	1	6	3	3	4		5	6	7	7	7	9	9	9	Output
0			3			6			9			12			15		18		

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

2	6	6	9	11	12	13	16	16	19	Next positions
0	1	2	3	4	5	6	7	8	9	

0	0		1	1	1	6	3	3	4	4	5	6	7	7	7	9	9	9	Output
0			3			6			9			12			15		18		

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

2	6	7	9	11	12	13	16	16	19	Next positions
0	1	2	3	4	5	6	7	8	9	

0	0		1	1	1	2	3	3	4	4	5	6	7	7	7	9	9	9	Output
0			3		6		9		12		15		18						

Distribution Counting Example (II)

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

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

0	3	6	7	9	11	12	13	16	16	Running sum of Counts
0	1	2	3	4	5	6	7	8	9	

3	6	7	9	11	12	13	16	16	19	Next positions
0	1	2	3	4	5	6	7	8	9	

0	0	0	1	1	1	2	3	3	4	4	5	6	7	7	7	9	9	9	Output
0			3			6			9		12			15			18		

Radix Sort

Idea: Sort keys *one character at a time*.

- We can use distribution counting for each digit.
- Can work either right to left (LSD radix sort) or left to right (MSD radix sort)
- LSD radix sort is venerable: used for punched cards. Example:

Initial: set, cat, cad, con, bat, can, be, let, bet

Pass 1 (by char #2)				Pass 2 (by char #1)				Pass 3 (by char #0)			
be	cad	can	con	bat	cat	can	cad	bet	cat	let	set
///	///	///	///	///	///	///	///	///	///	///	///
'b'	'd'	'n'	't'						'a'	'e'	'o'
be, cad, con, can, set, cat, bat, let, bet						cad, can, cat, bat, be, set, let, bet, con					
con bet cat be can bat cad let set								bat, be, bet, cad, can, cat, con, let, set			
'b' 'd' 'n' 't'						'a' 'e' 'o'					

MSD Radix Sort

- A bit more complicated: must keep lists from each step separate,
- But, we can stop processing 1-element lists

A	posn
* set, cat, cad, con, bat, can, be, let, bet	0
* bat, be, bet / cat, cad, con, can / let / set	1
bat / * be, bet / cat, cad, con, can / let / set	2
bat / be / bet / * cat, cad, con, can / let / set	1
bat / be / bet / * cat, cad, can / con / let / set	2
bat / be / bet / cad / can / cat / con / let / set	

- Here, slashes divide partially sorted sublists, which will never be moved into the spaces occupied by other sublists.
- Asterisks mark a sublist to be sorted on some character position.

Performance of Radix Sort

- Radix sort takes $\Theta(B)$ time where B is *the total size of the key data (total characters)*.
- Up until now, we have measured sort times as functions of #records.
- How to compare these two different measures?
- To have N different records, one must have keys at least $\Theta(\lg N)$ long [why?]
- Furthermore, comparisons actually take time $\Theta(K)$ where K is the size of key in worst case [why?]
- So $N \lg N$ comparisons really means $N(\lg N)^2$ operations.
- While radix sort would take $B = N \lg N$ time for N records with minimal-length ($\lg N$) keys.
- On the other hand, we must work to get good constant factors with radix sort.

And Don't Forget Search Trees

Idea: A search tree is in sorted order, when read in inorder.

- Need *balance* to really use for sorting [next topic].
- Given balance, same performance as heapsort: N insertions in time $\lg N$ each, plus $\Theta(N)$ to traverse, gives

$$\Theta(N + N \lg N) = \Theta(N \lg N)$$

Summary

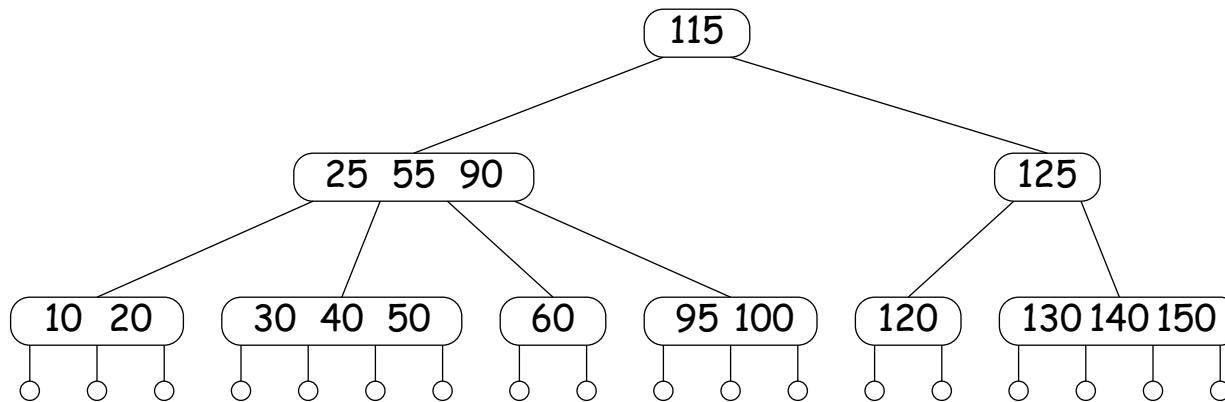
- Insertion sort: $\Theta(Nk)$ comparisons and moves, where k is maximum amount data is displaced from final position.
 - Good for small datasets or almost ordered data sets.
- Quicksort: $\Theta(N \lg N)$ with good constant factor if data is not pathological. Worst case $O(N^2)$.
- Merge sort: $\Theta(N \lg N)$ guaranteed. Good for external sorting.
- Heapsort, treesort with guaranteed balance: $\Theta(N \lg N)$ guaranteed.
- Radix sort, distribution sort: $\Theta(B)$ (number of bytes). Also good for external sorting.

CS61B Lecture #29

Balanced Search: The Problem

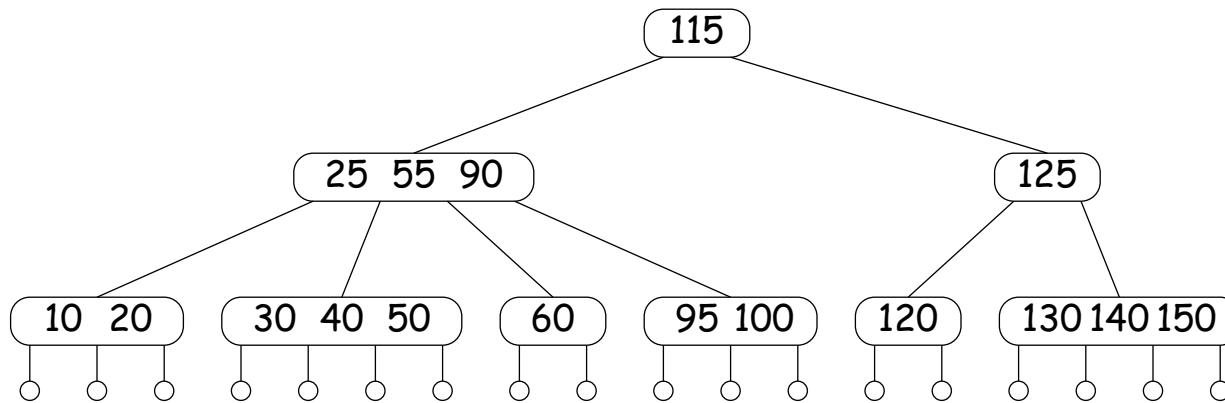
- Why are search trees important?
 - Insertion/deletion is fast (on every operation, unlike hash table, which has to expand from time to time).
 - They support range queries and sorting (unlike hash tables).
- But $O(\lg N)$ performance from binary search tree requires keys to be roughly divided by some constant > 1 at each inner node.
- In other words, the tree must be “bushy”
- “Stringy” trees (where most inner nodes have one child) perform like linked lists.
- It suffices that the heights of the subtrees of nodes that are not near the bottom always differ in height by no more than some constant factor $C > 0$.

Example of Direct Approach: B-Trees



- An *order M B-tree* is an M -ary search tree, $M > 2$, where non-root nodes have at least $\lceil M/2 \rceil$ children.
- For example, if $M = 4$, non-root nodes have at least 2 nodes, so we also call these $(2, 4)$ (or 2-4) trees.
- Obeys the search-tree property:
 - Keys are sorted in each node.
 - Each key separates two subtrees with all keys in the left subtree of a key, K , being $< K$, and all to the right being $> K$.
- Children at the bottom of tree are all empty (don't really exist) and equidistant from the root.
- Searching is a simple generalization of binary search.

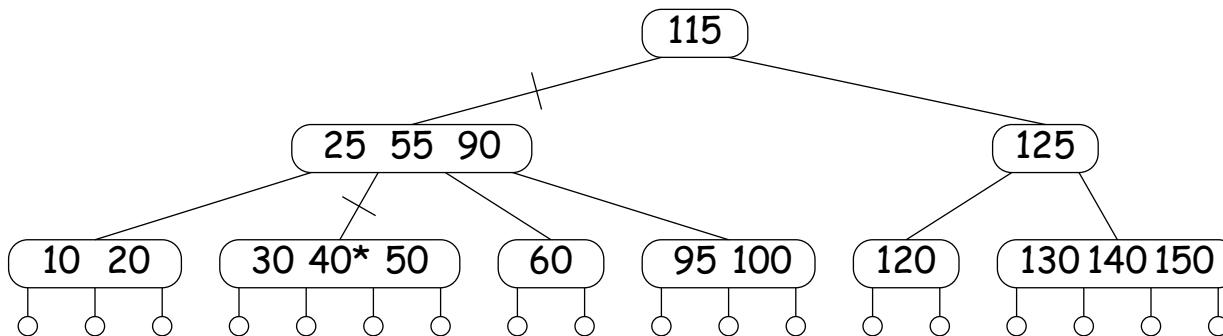
Example of Direct Approach: B-Trees



Idea: If the tree grows/shrinks only at the root, then the two sides always have same height.

- Each node, except the root, has from $\lceil M/2 \rceil$ to M children, and one key “between” each two children.
- The root has from 2 to M children (in a non-empty tree).
- Insertion: add to nodes just above the (empty) leaves; split overfull nodes as needed, moving one key up to its parent.

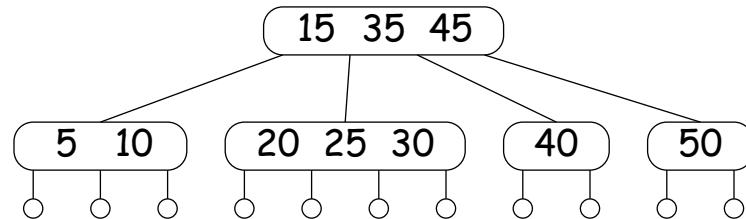
Sample Order 4 B-tree ((2,4) Tree)



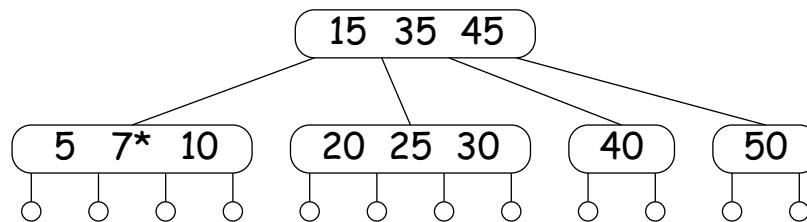
- Crossed lines show the path for finding 40.
- Keys on either side of each child pointer in the path bracket 40.
- Each node has at least 2 children, and all leaves (little circles) are at the bottom, so the height must be $O(\lg N)$.
- In real life, B-trees are stored on secondary storage, and the order is typically much bigger: comparable to the size of a disk sector, page, or other convenient unit of I/O.

Inserting in (2, 4) tree (Simple Case)

- Start:

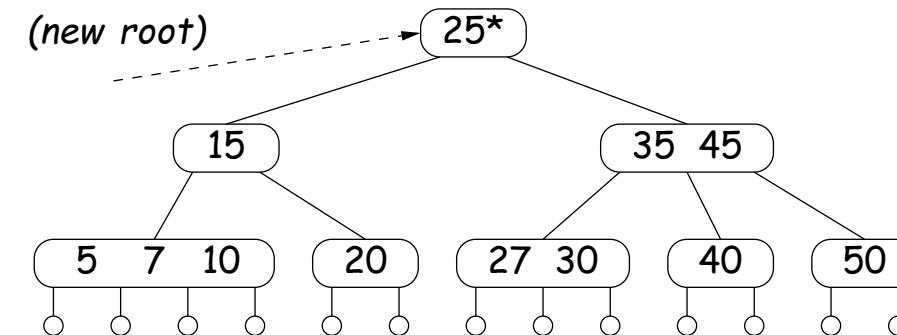
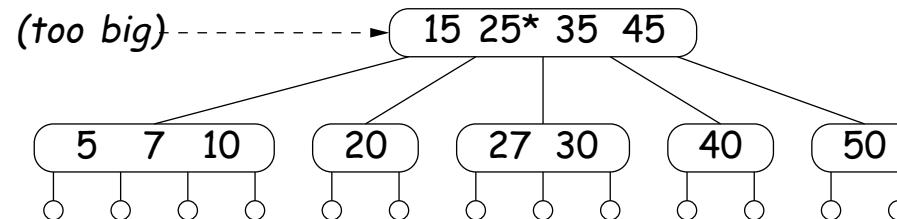
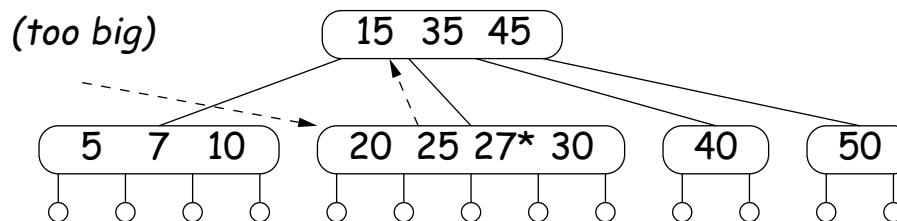


- Insert 7:



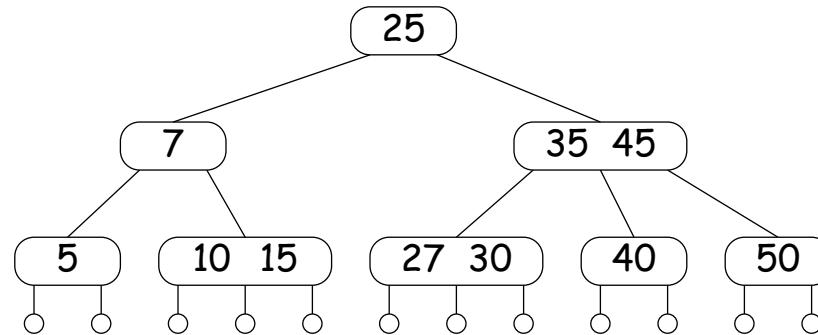
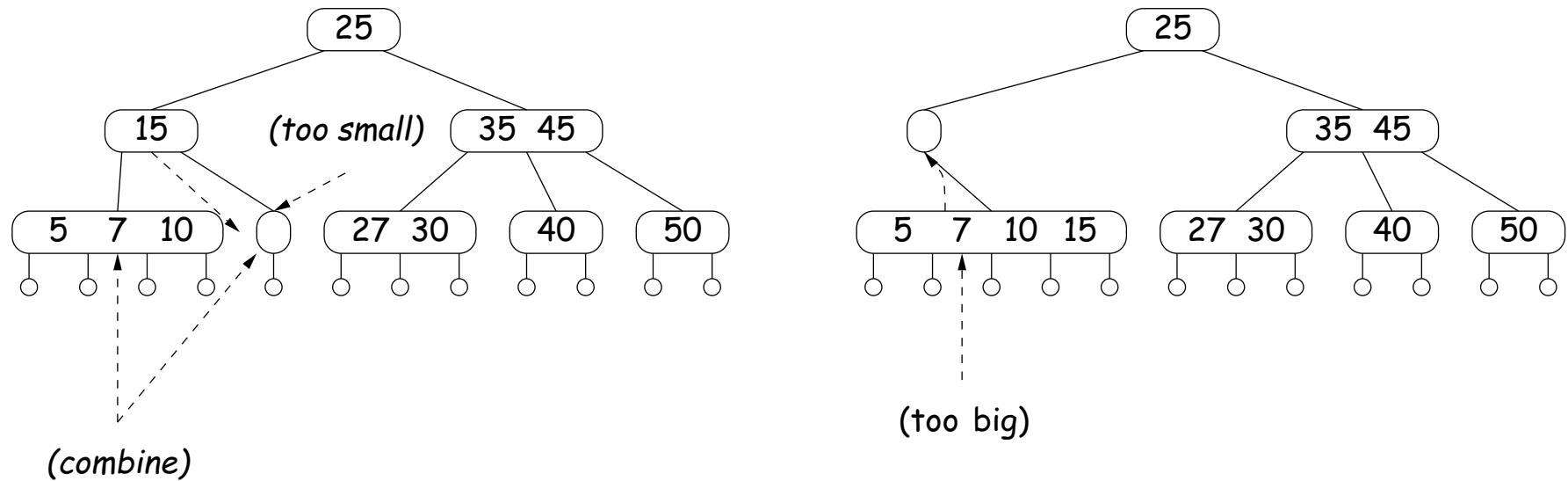
Inserting in B-Tree (Splitting)

- Insert 27:



Deleting Keys from B-tree

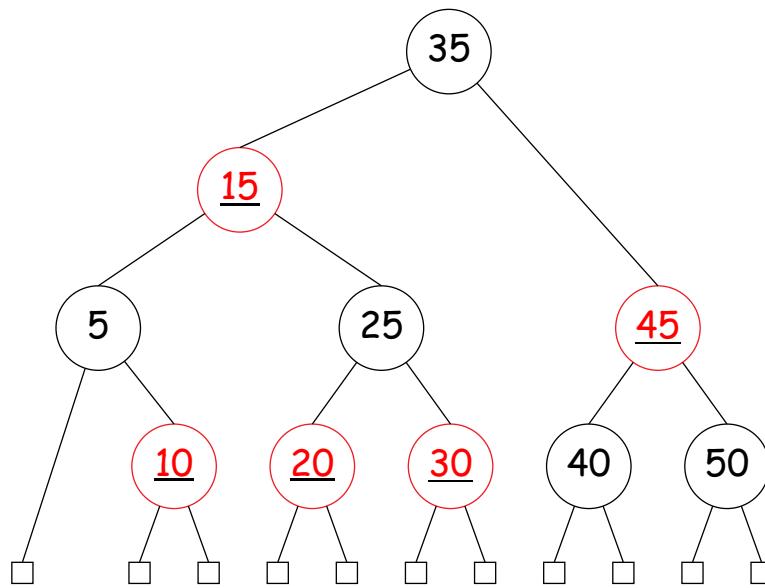
- Remove 20 from last tree.



Red-Black Trees

- A *Red-black tree* is a binary search tree with additional constraints that limit how unbalanced it can be.
- Thus, searching is always $O(\lg N)$.
- Used for Java's TreeSet and TreeMap types.
- When items are inserted or deleted, the tree is *rotated* and *recolored* as needed to restore balance.

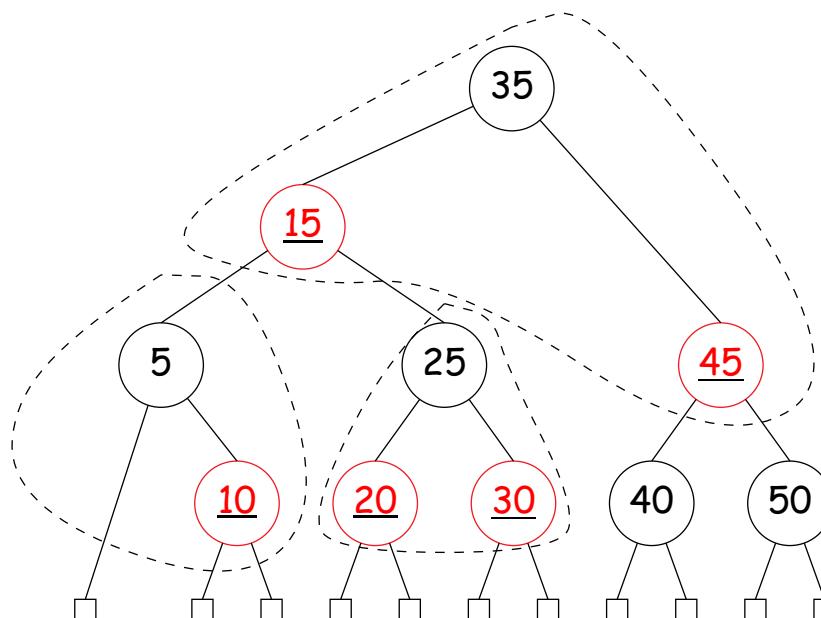
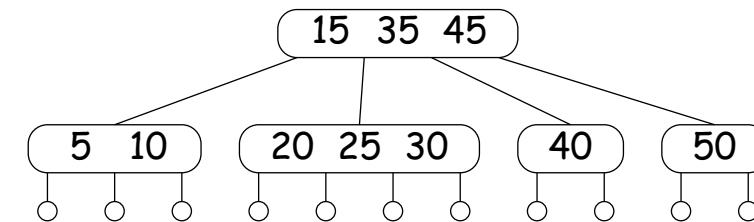
Red-Black Tree Constraints



1. Each node is (conceptually) colored red or black.
 2. Root is black.
 3. Every leaf node contains no data (as for B-trees) and is black.
 4. Every leaf has same number of black ancestors.
 5. Every internal node has two children.
 6. Every red node has two black children.
- Conditions 4, 5, and 6 guarantee $O(\lg N)$ searches.

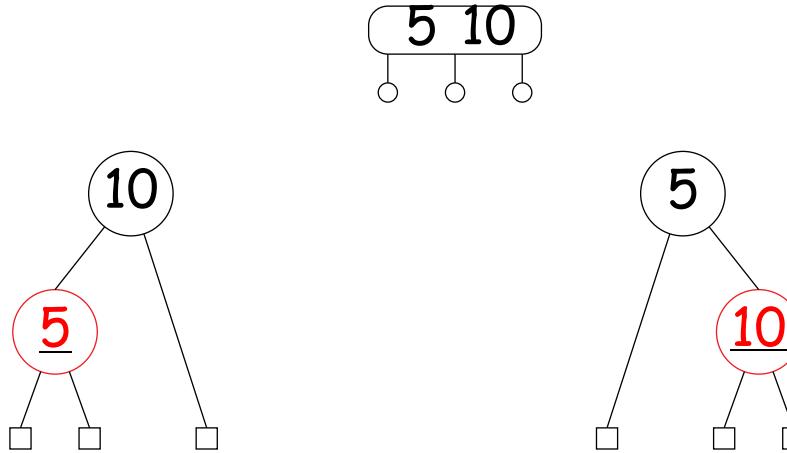
Red-Black Trees and (2,4) Trees

- Every red-black tree corresponds to a (2,4) tree, and the operations on one correspond to those on the other.
- Each node of (2,4) tree corresponds to a cluster of 1-3 red-black nodes in which the top node is black and any others are red.



Additional Constraints: Left-Leaning Red-Black Trees

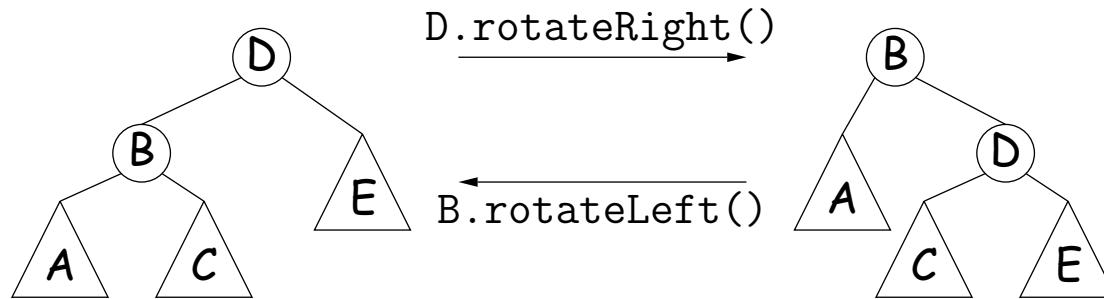
- A node in a (2,4) or (2,3) tree with three children may be represented in two different ways in a red-black tree:



- We can considerably simplify insertion and deletion in a red-black tree by always choosing the option on the left.
- With this constraint, there is a one-to-one relationship between (2,4) trees and red-black trees.
- The resulting trees are called *left-leaning red-black trees*.
- As a further simplification, let's restrict ourselves to red-black trees that correspond to (2,3) trees (whose nodes have no more than 3 children), so that no red-black node has two red children.

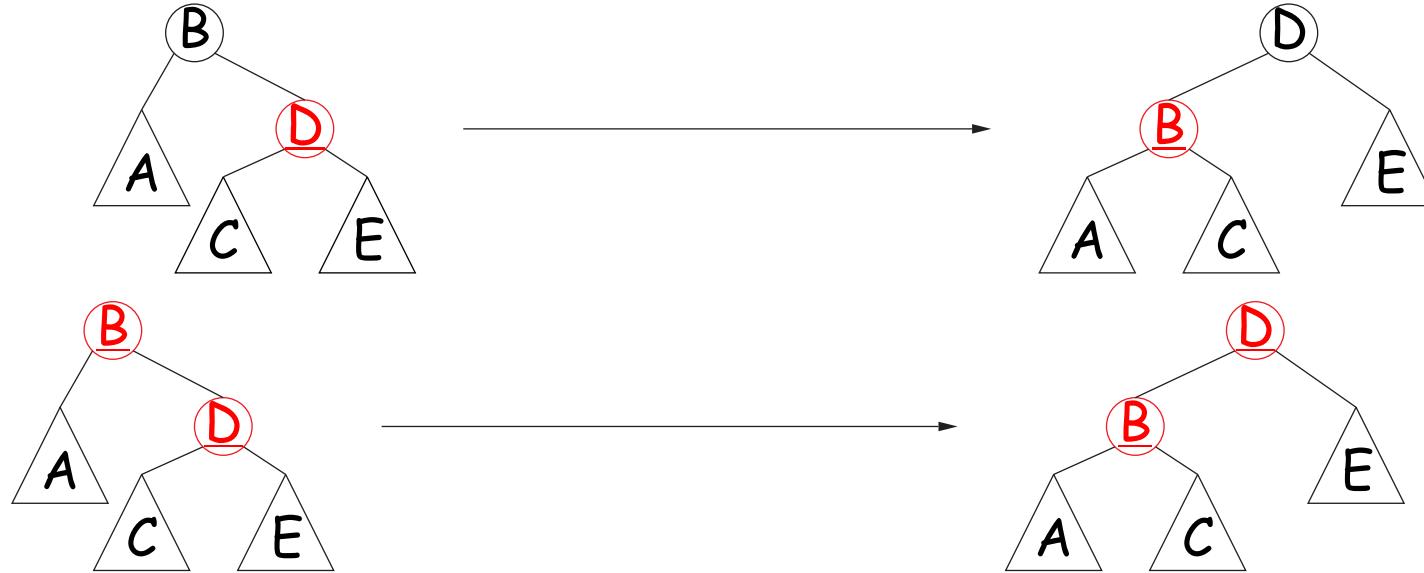
Red-Black Insertion and Rotations

- Insert at the bottom just as for binary trees (color the insertion red except when the tree is initially empty).
- Then rotate (and recolor) to restore the red-black property, and thus rebalance.
- Rotation of trees *preserves* the binary tree property, but changes balance.



Rotations and Recolorings

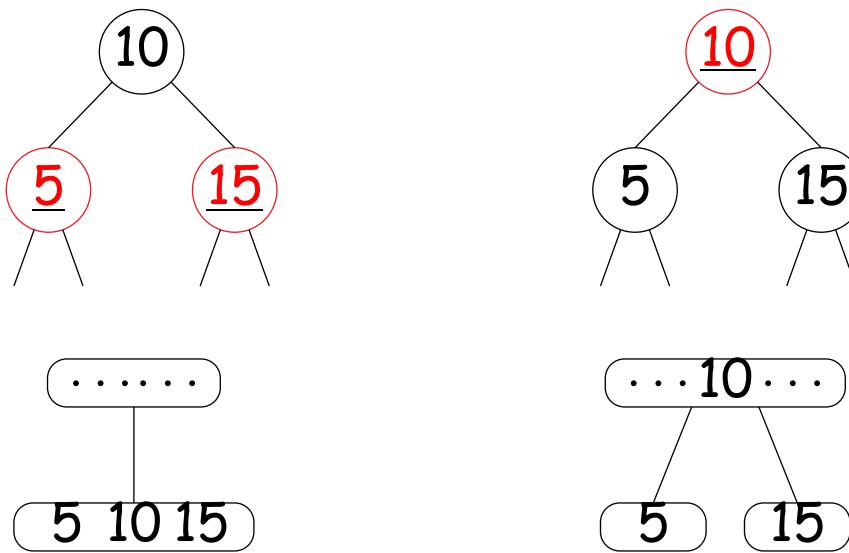
- For our purposes, we'll augment the general rotation algorithms with some recoloring.
- Transfer the color from the original root to the new root, and color the original root red. Examples:



- Neither of these changes the number of black nodes along any path between the root and the leaves.

Splitting by Recoloring

- Our algorithms will temporarily create nodes with too many children, and then split them up.
- A simple recoloring allows us to split nodes. We'll call it `colorFlip`:



- Here, key 10 joins the parent node, splitting the original.

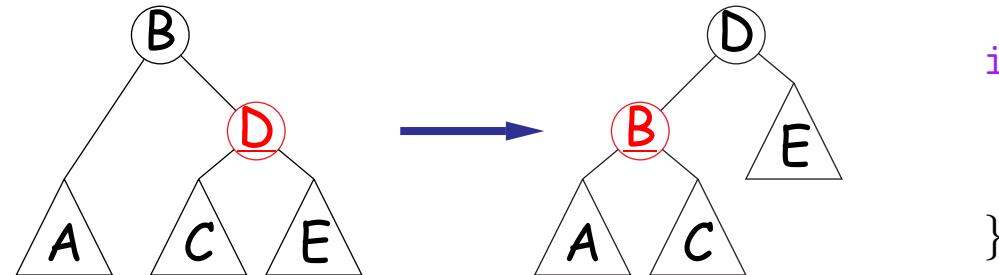
The Algorithm (Sedgewick)

- We posit a binary-tree type RBTree: basically ordinary BST nodes plus color.
- Insertion is the same as for ordinary BSTs, but we add some fixups to restore the red-black properties.

```
RBTree insert(RBTree tree, KeyType key) {  
    if (tree == null)  
        return new RBTree(key, null, null, RED);  
    int cmp = key.compareTo(tree.label());  
    else if (cmp < 0) tree.setLeft(insert(tree.left(), key));  
    else                  tree.setRight(insert(tree.right(), key));  
  
    return fixup(tree);      // Only line that's all new!  
}
```

Fixing Up the Tree

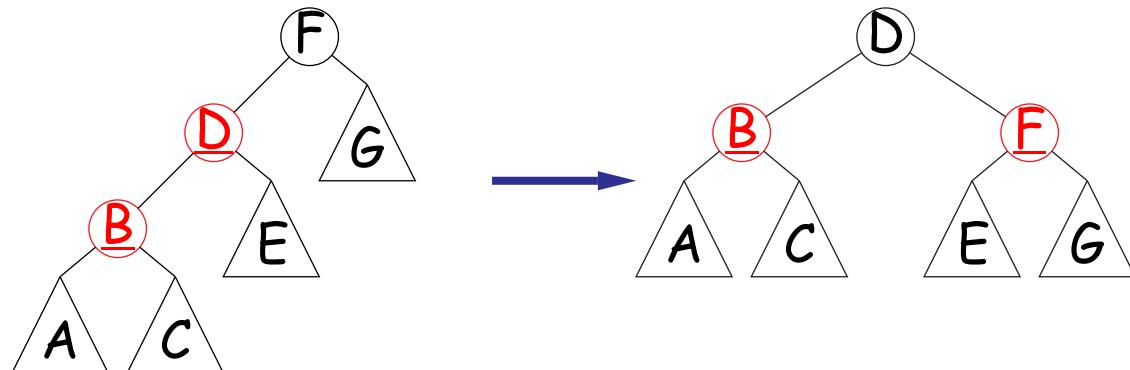
- As we return back up the BST, we restore the left-leaning red-black properties, and limit ourselves to red-black trees that correspond to (2,3) trees by applying the following (in order) to each node:
- Fixup 1: Convert right-leaning trees to left-leaning:



```
if (tree.right().isRed()
    && tree.left().isBlack()) {
    tree.rotateLeft();
}
```

Sometimes, node B will be red, so that both B and D end up red. This is fixed by...

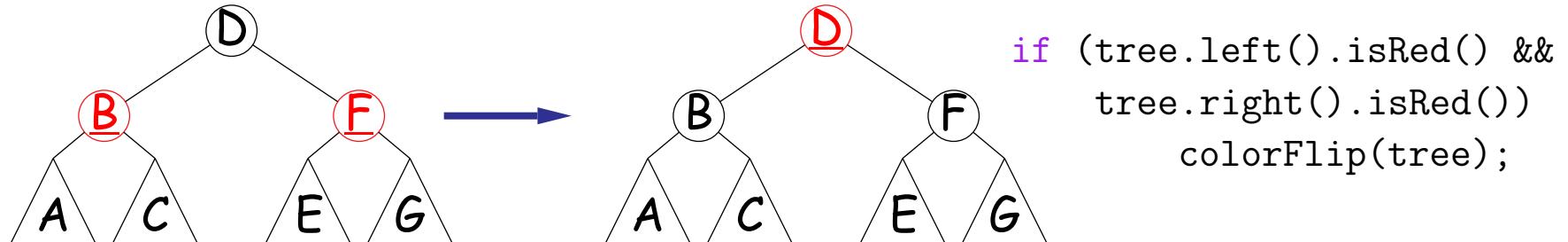
- Fixup 2: Rotate linked red nodes into a normal 4-node (temporarily).



```
if (tree.left().isRed() &&
    tree.left().left().isRed())
    tree.rotateRight();
```

Fixing Up the Tree (II)

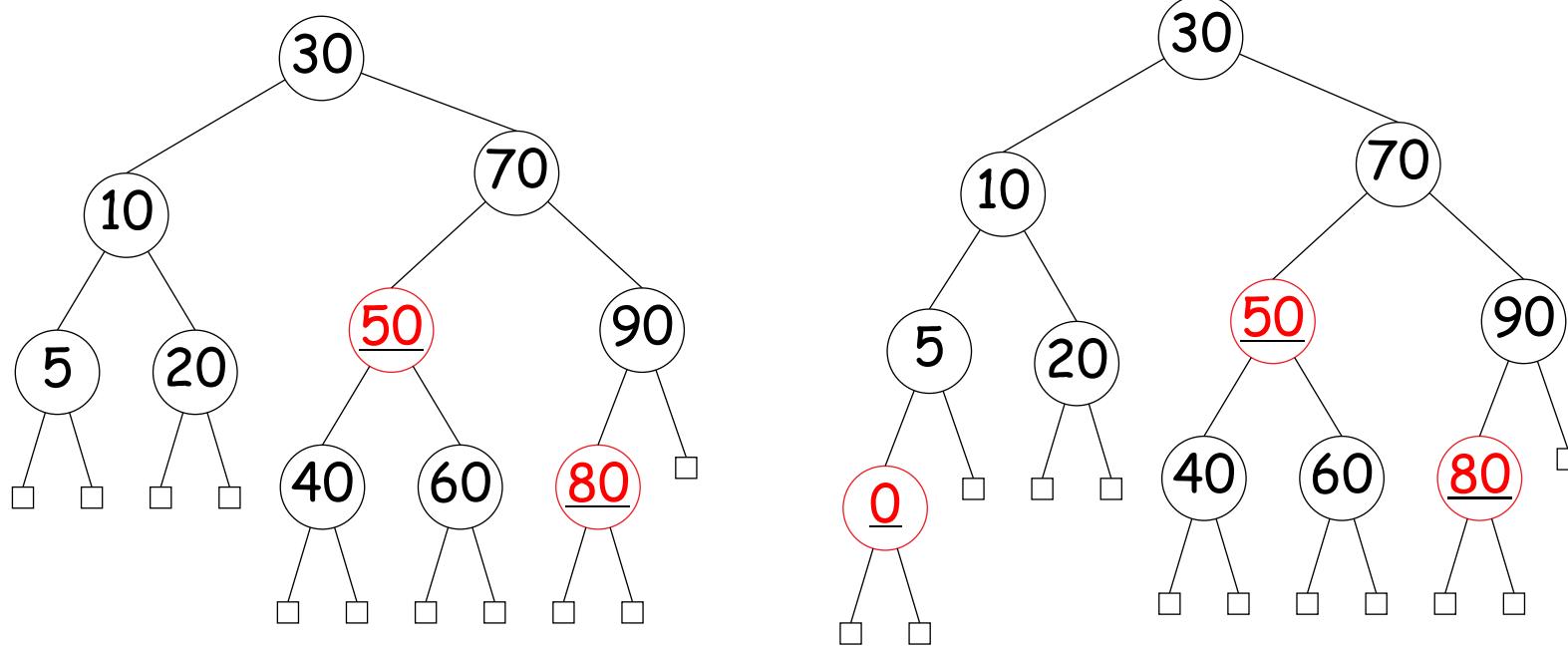
- Fixup 3: Break up 4-nodes into 3-nodes or 2-nodes.



- Fixup 4: As a result of other fixups, or of insertion into the empty tree, the root may end up red, so color the root black after the rest of insertion and fixups are finished. (Not part of the fixup function; just done at the end).

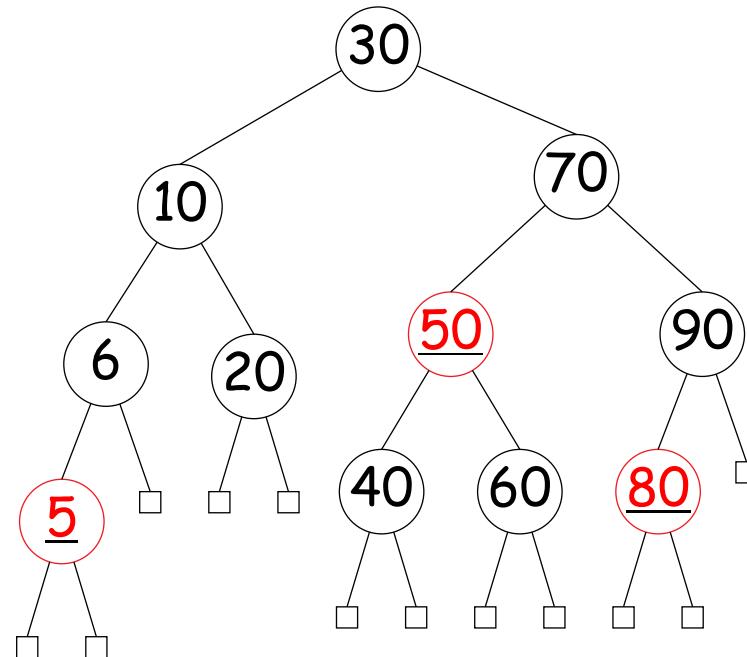
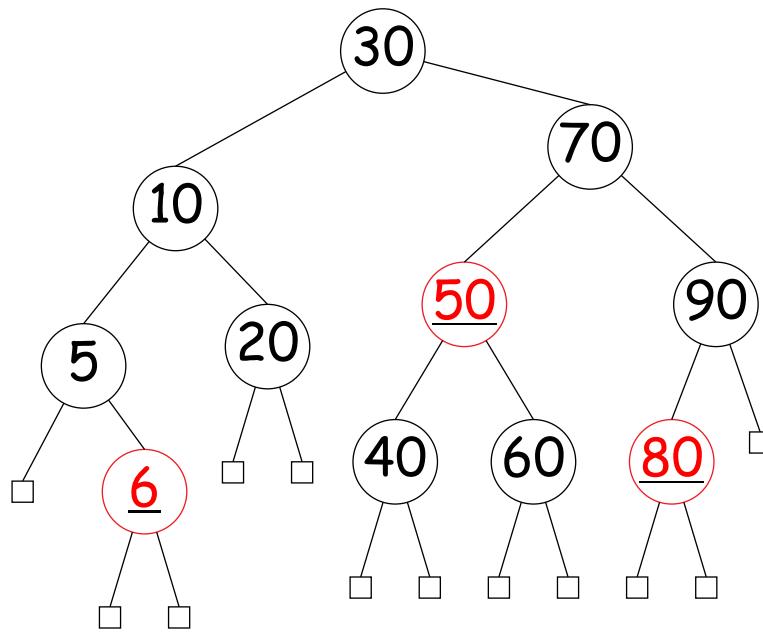
Example of Left-Leaning 2-3 Red-Black Insertion

- Insert 0 into initial tree on the left. No fixups needed.



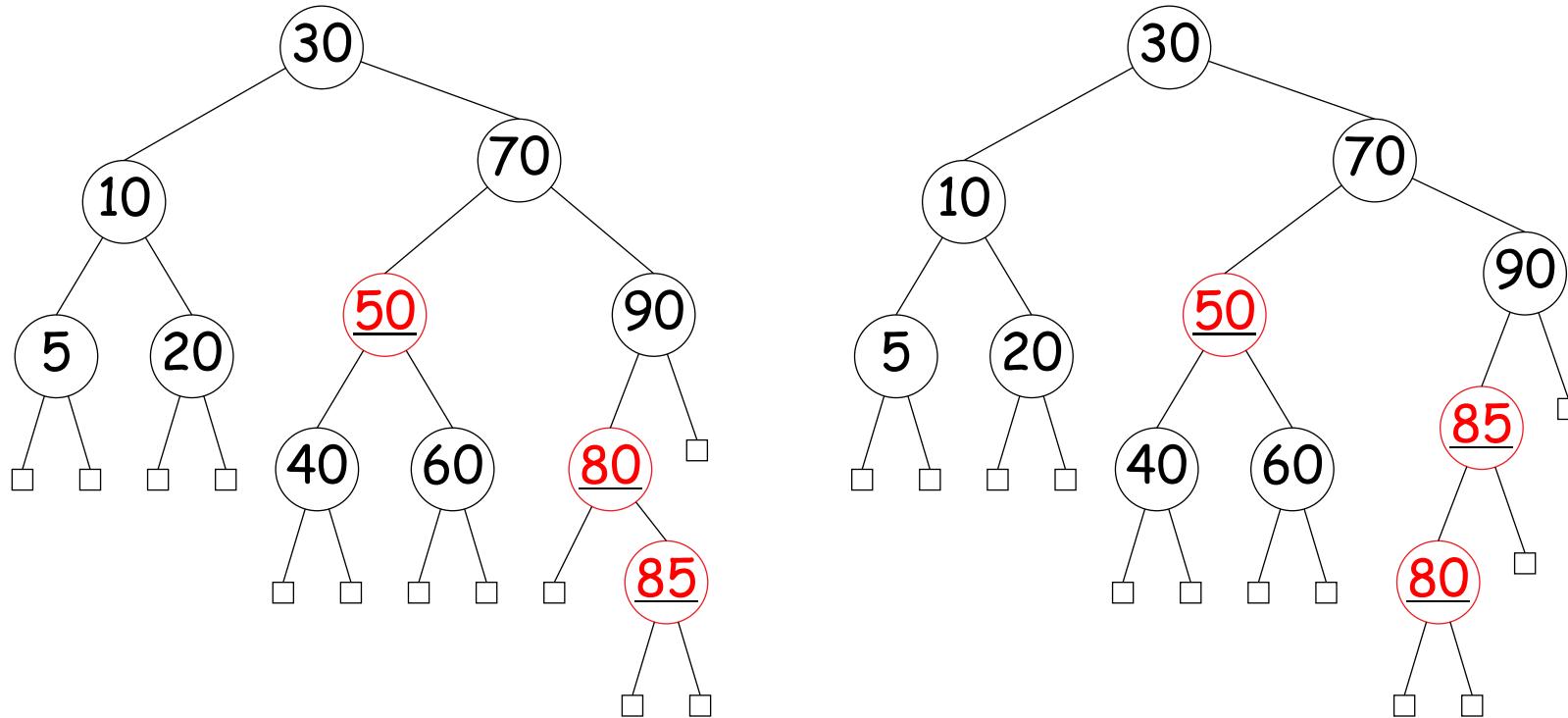
Insertion Example (II)

- Instead of 0, let's insert 6, leading to the tree on the left. This is right-leaning, so apply Fixup 1:



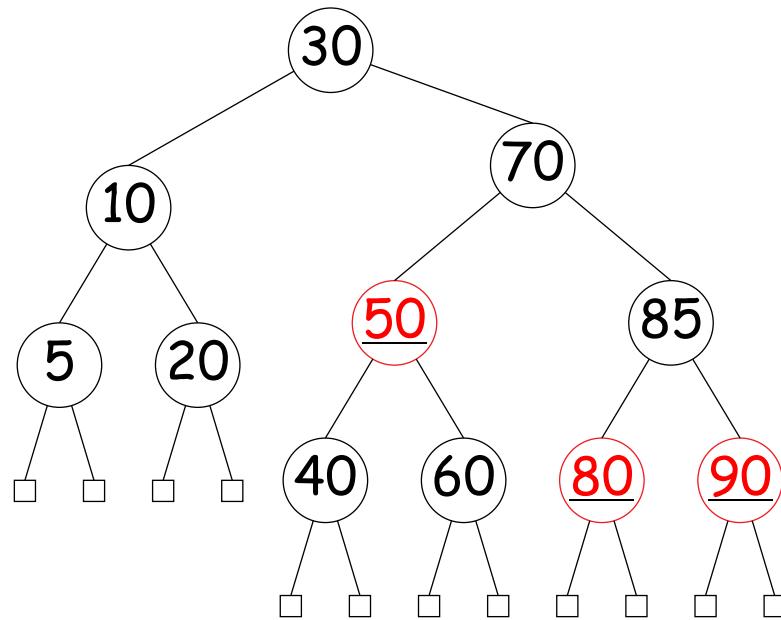
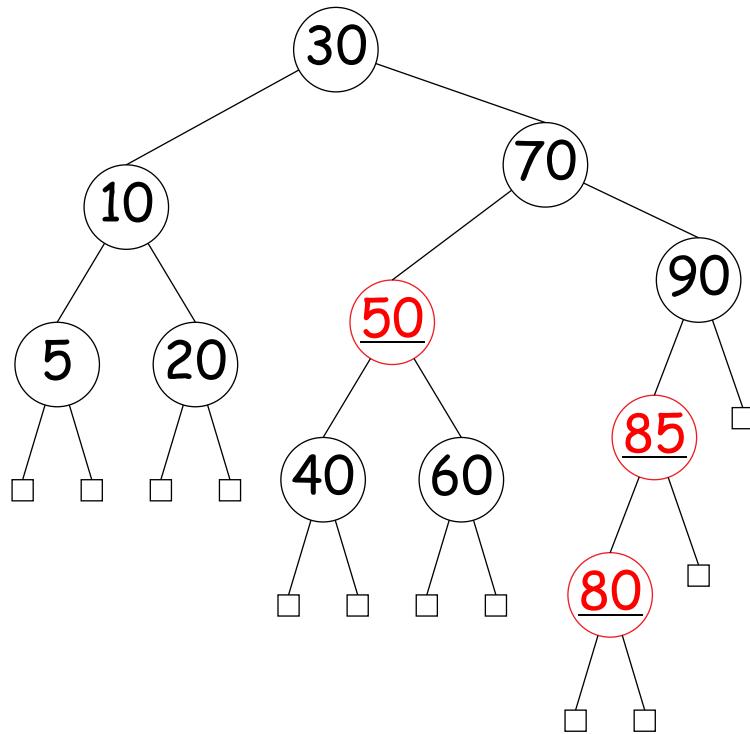
Insertion Example (III)

- Now consider inserting 85. We need fixup 1 first.



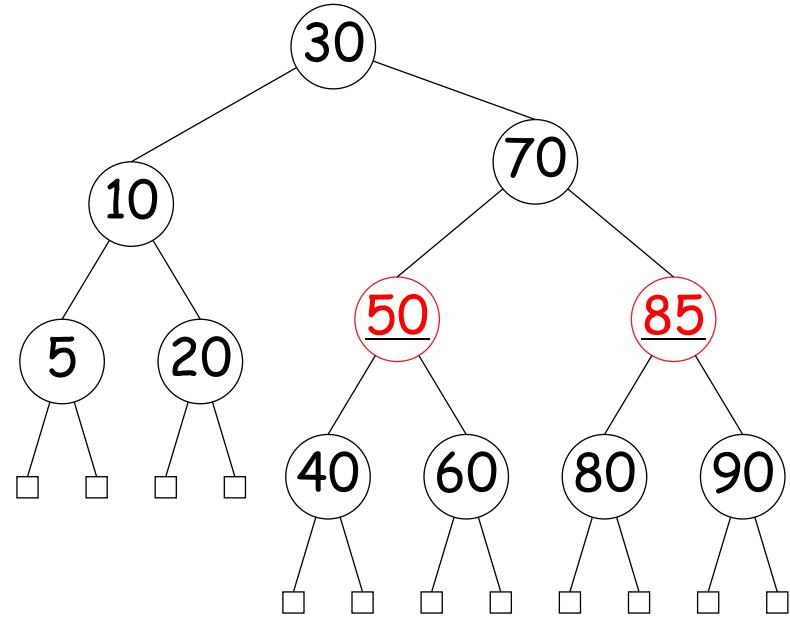
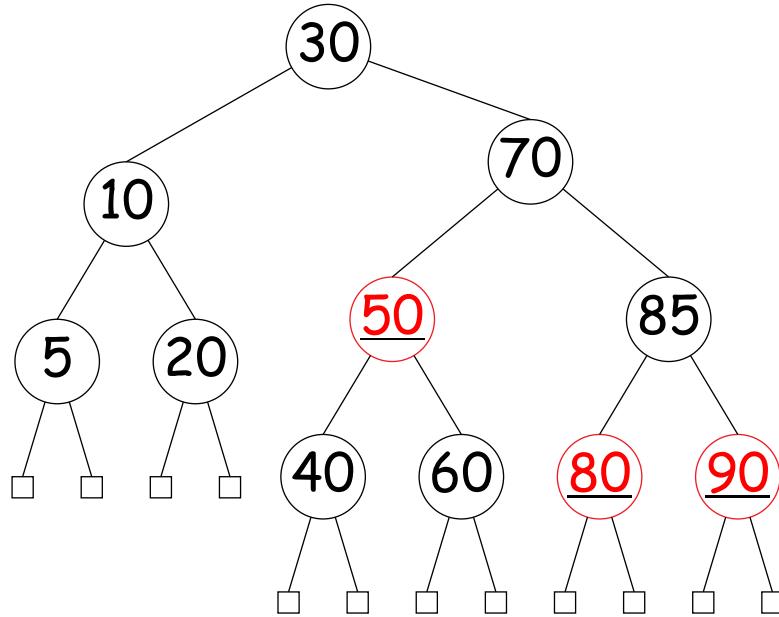
Insertion Example (IIIa)

- Now apply fixup 2.



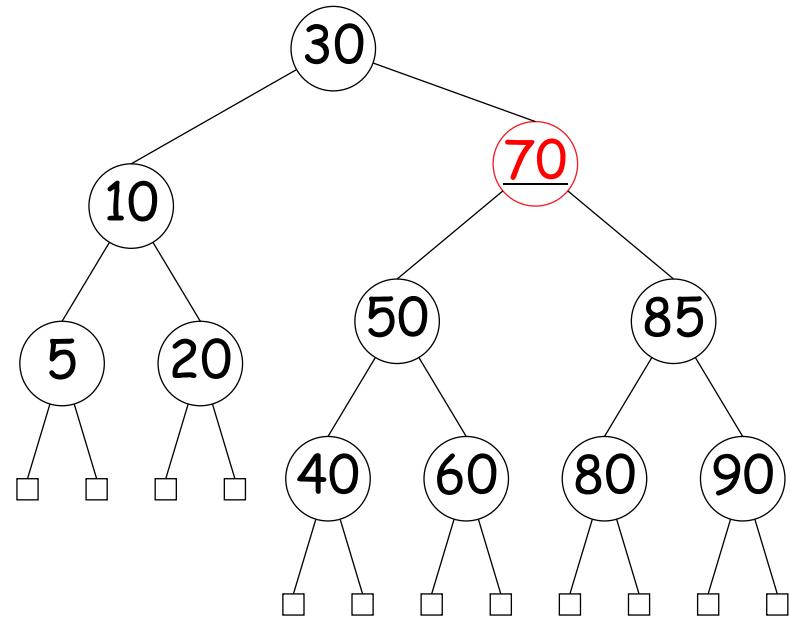
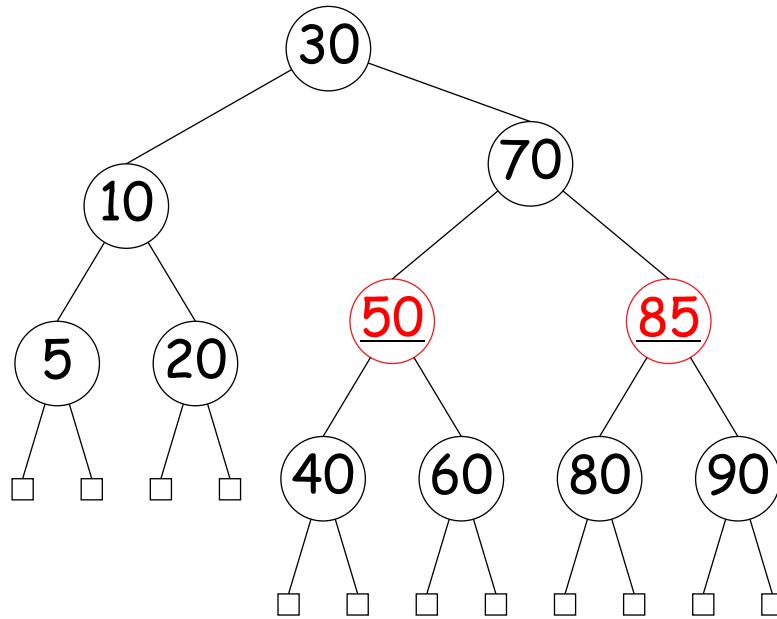
Insertion Example (IIIb)

- This gives us a 4-node, so apply fixup 3.



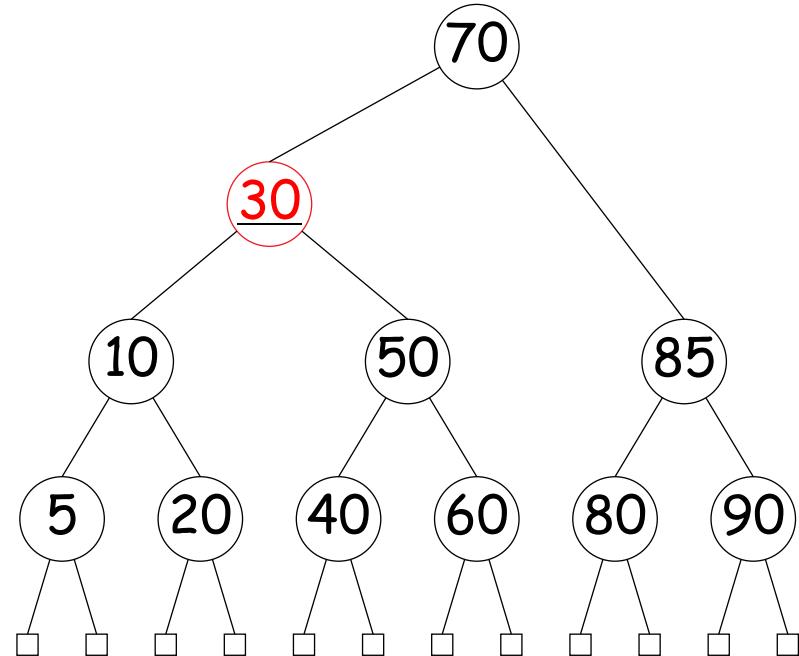
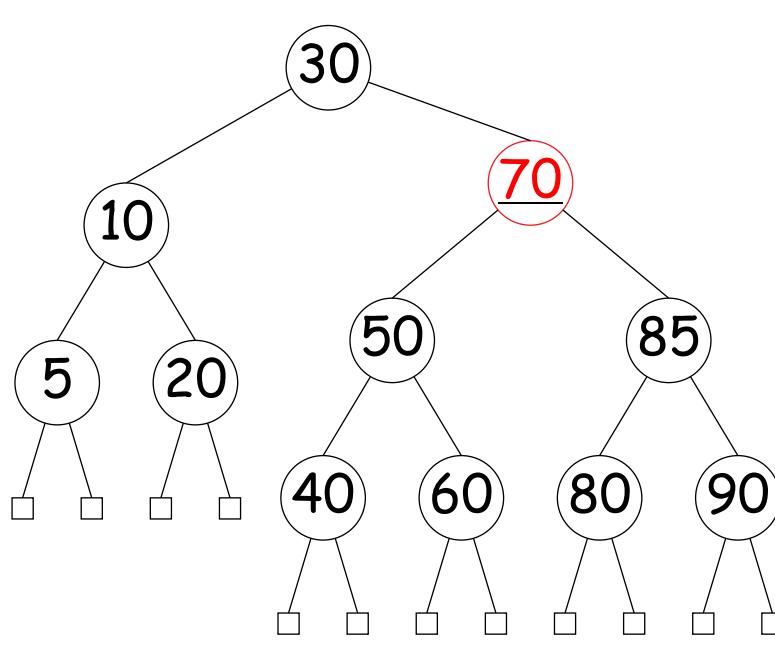
Insertion Example (IIIc)

- This gives us another 4-node, so apply fixup 3 again.



Insertion Example (III^d)

- This gives us a right-leaning tree, so apply fixup 1.



CS61B Lecture #31

Today:

- More balanced search structures (*DS(IJ)*, Chapter 9)

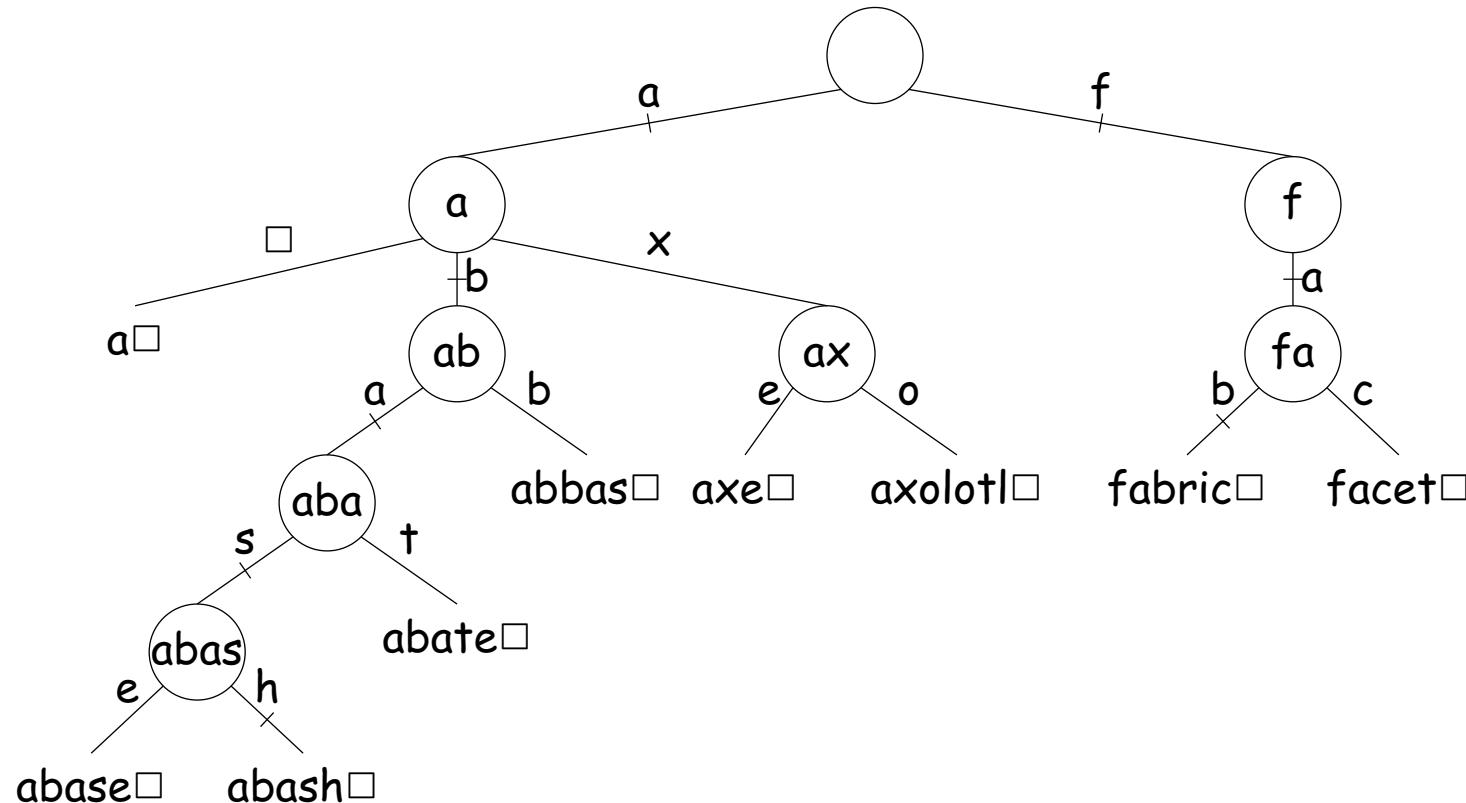
Really Efficient Use of Keys: the Trie

- We haven't said much about the cost of comparisons, generally treating the cost as constant.
- For strings, the worst case is the length of string.
- Therefore we should throw an extra factor of the key length, L , into costs:
 - $\Theta(M)$ comparisons really means $\Theta(ML)$ operations.
 - So to look for key X , we keep looking at same chars of X M times.
- Can we do better? Can we get search cost to be $O(L)$?

Idea: Make a *multi-way decision tree*, with one decision per character of key.

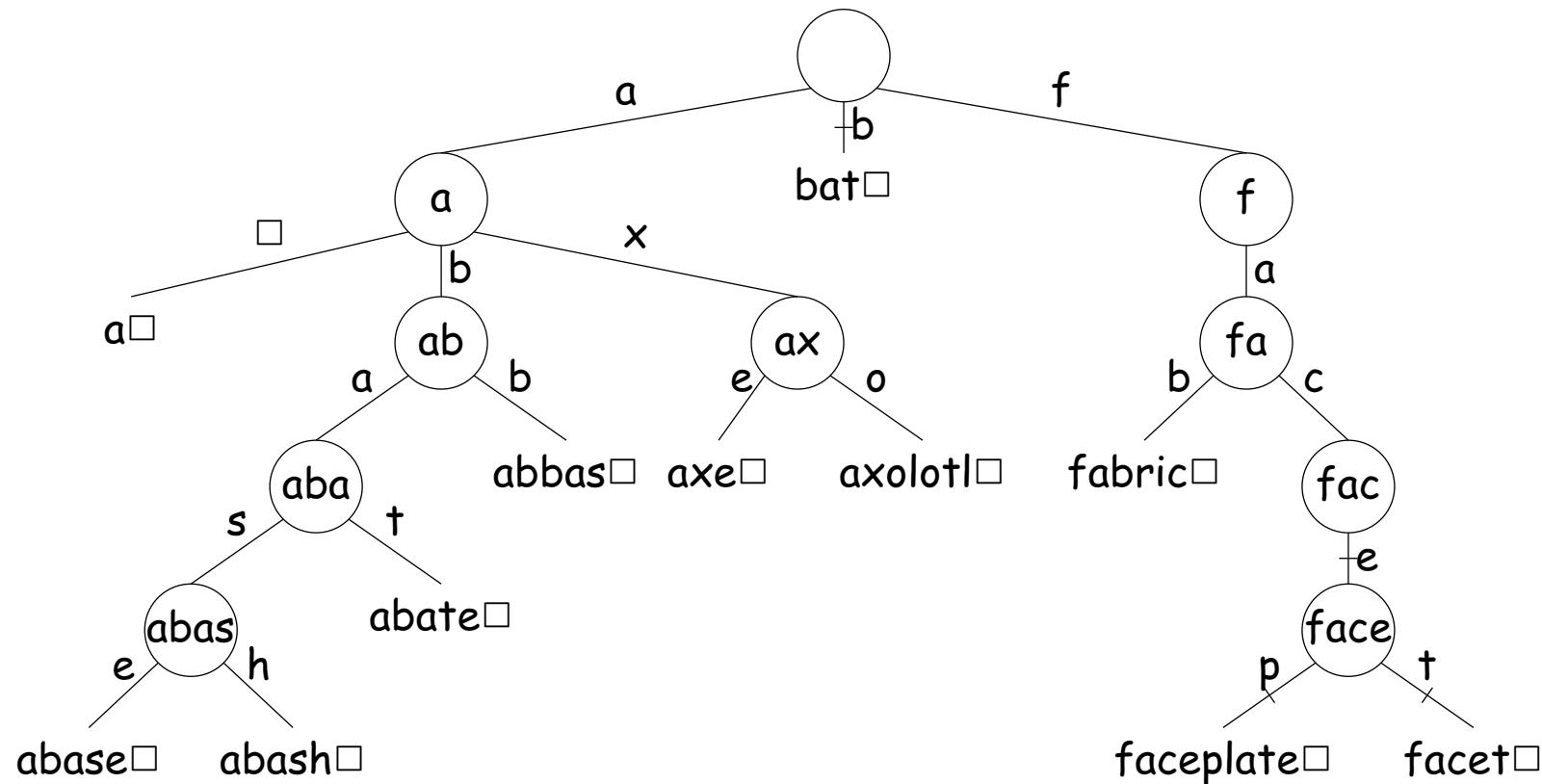
The Trie: Example

- Set of keys
 - {a, abase, abash, abate, abbas, axolotl, axe, fabric, facet}
- Ticked lines show paths followed for "abash" and "fabric"
- Each internal node corresponds to a possible prefix.
- Characters in path to node = that prefix.



Adding Item to a Trie

- Result of adding bat and faceplate.
- New edges ticked.



A Side-Trip: Scrunching

- For speed, obvious implementation for internal nodes is array indexed by character.
- Gives $O(L)$ performance, L length of search key.
- [Looks as if independent of N , number of keys. Is there a dependence?]
- Problem: arrays are *sparsely populated* by non-null values—waste of space.

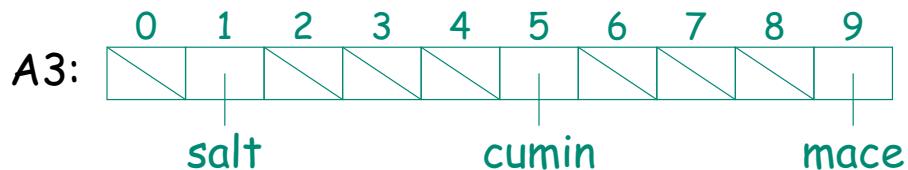
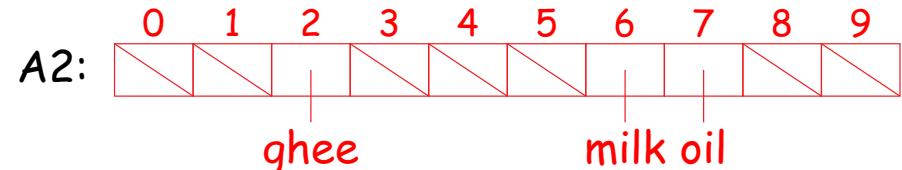
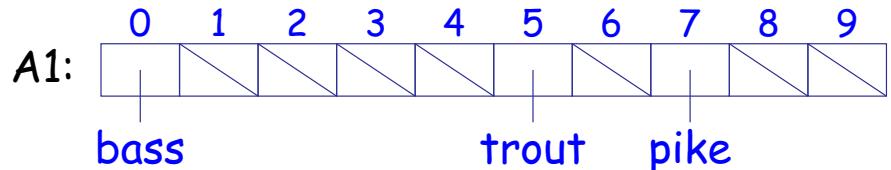
Idea: Put the arrays on top of each other!

- Use null (0, empty) entries of one array to hold non-null elements of another.
- Use extra markers to tell which entries belong to which array.

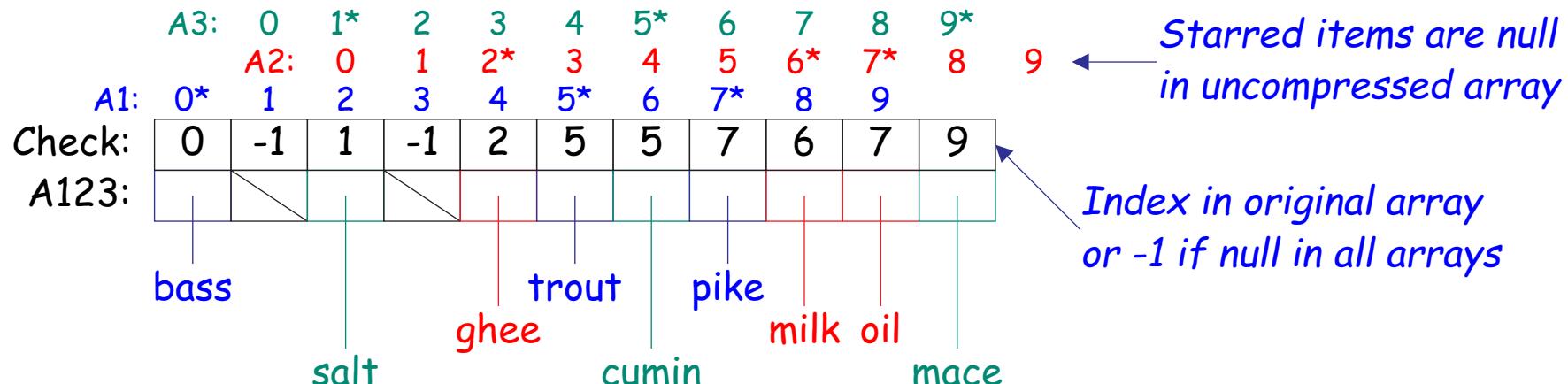
Scrunching Example

Small example: (unrelated to Tries on preceding slides)

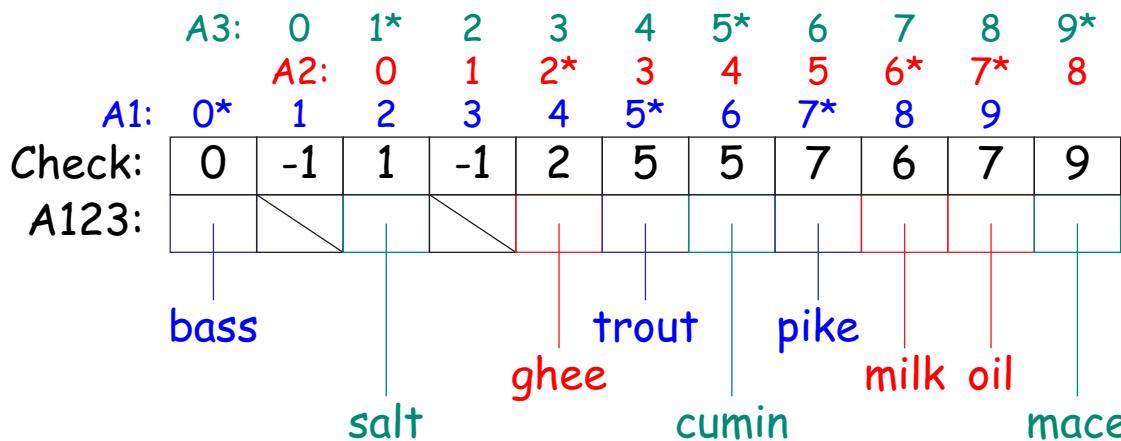
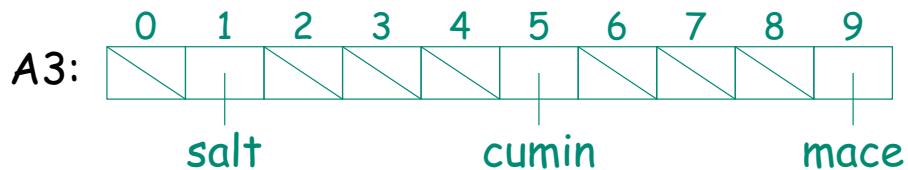
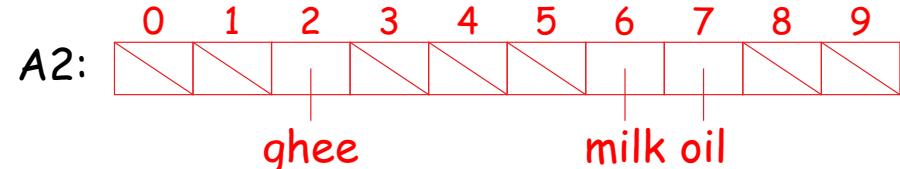
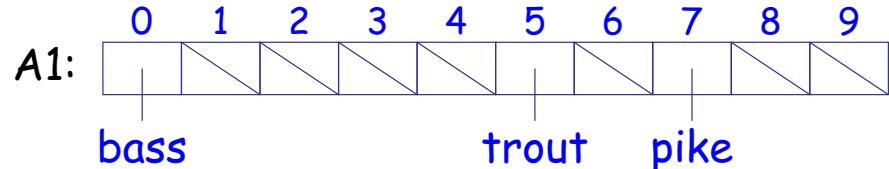
- Three arrays, each indexed 0..9



- Now overlay them, but keep track of the original index of each item:



Scrunching Example (contd.)



```
/* A1[i] == */ (Check[i] == i) ? A123[i] : null;
/* A2[i] == */ (Check[i + 2] == i) ? A123[i + 2] : null;
/* A3[i] == */ (Check[i + 1] == i) ? A123[i + 1] : null;
```

Starred items are null
in uncompressed array

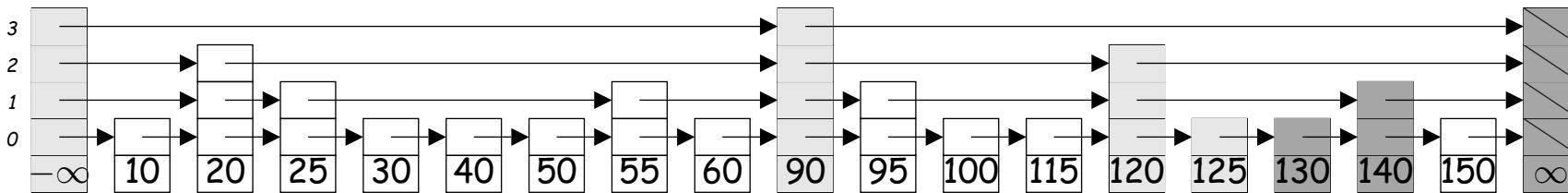
Index in original array
or -1 if null in all arrays

Practicum

- The scrunching idea is cute, but
 - Not so good if we want to expand our trie.
 - A bit complicated.
 - Actually more useful for representing large, sparse, fixed tables with many rows and columns.
- Furthermore, number of children in trie tends to drop drastically when one gets a few levels down from the root.
- So in practice, might as well use linked lists to represent set of node's children...
- ...but use arrays for the first few levels, which are likely to have more children.

Probabilistic Balancing: Skip Lists

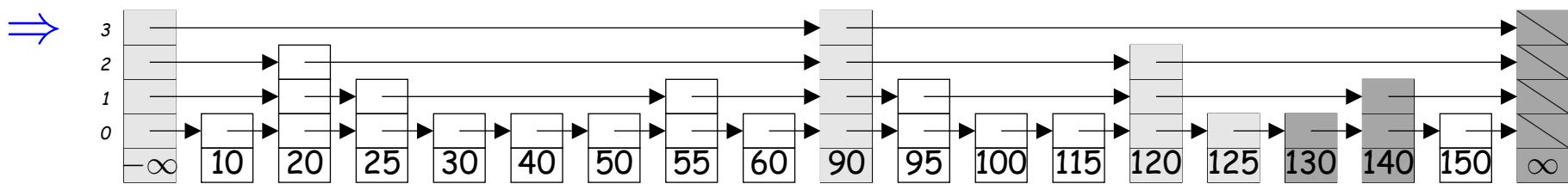
- A *skip list* can be thought of as a kind of n-ary search tree in which we choose to put the keys at "random" heights.
- More often thought of as an ordered list in which one can skip large segments.
- Typical example:



- To search, start at top layer on left, search until next step would overshoot, then go down one layer and repeat.
- In list above, we search for 125 and 127. Gray nodes are looked at; darker gray nodes are overshoots.
- Heights of the nodes were chosen randomly so that there are about $1/2$ as many nodes that are $> k$ high as there are that are k high.
- Makes searches fast *with high probability*.

Probabilistic Balancing: Skip Lists

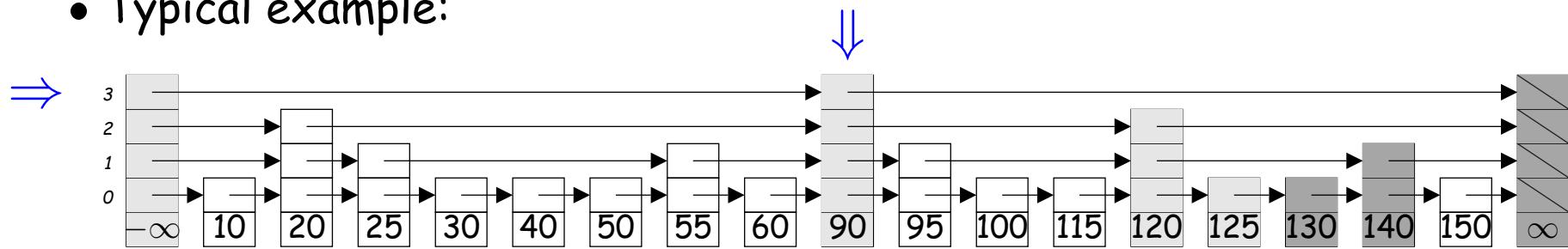
- A *skip list* can be thought of as a kind of n-ary search tree in which we choose to put the keys at "random" heights.
- More often thought of as an ordered list in which one can skip large segments.
- Typical example:



- To search, start at top layer on left, search until next step would overshoot, then go down one layer and repeat.
- In list above, we search for 125 and 127. Gray nodes are looked at; darker gray nodes are overshoots.
- Heights of the nodes were chosen randomly so that there are about $1/2$ as many nodes that are $> k$ high as there are that are k high.
- Makes searches fast *with high probability*.

Probabilistic Balancing: Skip Lists

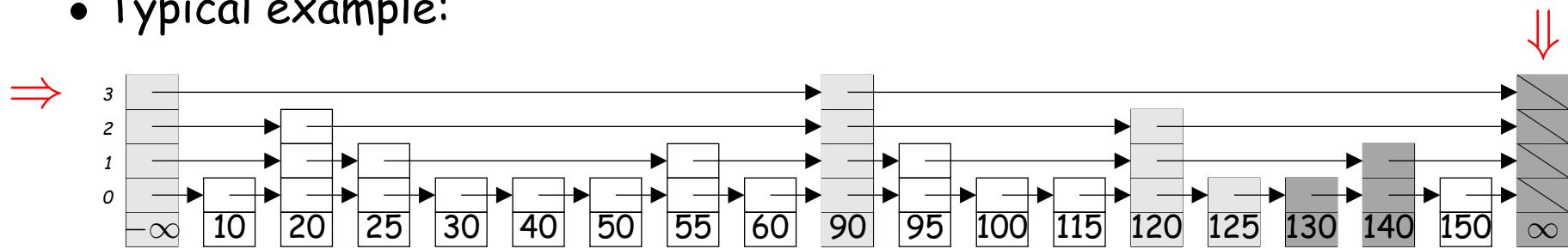
- A *skip list* can be thought of as a kind of n-ary search tree in which we choose to put the keys at "random" heights.
- More often thought of as an ordered list in which one can skip large segments.
- Typical example:



- To search, start at top layer on left, search until next step would overshoot, then go down one layer and repeat.
- In list above, we search for 125 and 127. Gray nodes are looked at; darker gray nodes are overshoots.
- Heights of the nodes were chosen randomly so that there are about $1/2$ as many nodes that are $> k$ high as there are that are k high.
- Makes searches fast *with high probability*.

Probabilistic Balancing: Skip Lists

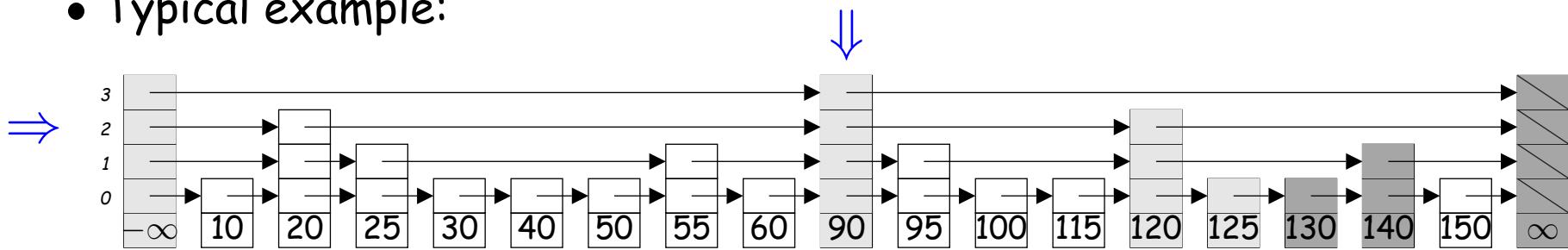
- A *skip list* can be thought of as a kind of n-ary search tree in which we choose to put the keys at "random" heights.
- More often thought of as an ordered list in which one can skip large segments.
- Typical example:



- To search, start at top layer on left, search until next step would overshoot, then go down one layer and repeat.
- In list above, we search for 125 and 127. Gray nodes are looked at; darker gray nodes are overshoots.
- Heights of the nodes were chosen randomly so that there are about $1/2$ as many nodes that are $> k$ high as there are that are k high.
- Makes searches fast *with high probability*.

Probabilistic Balancing: Skip Lists

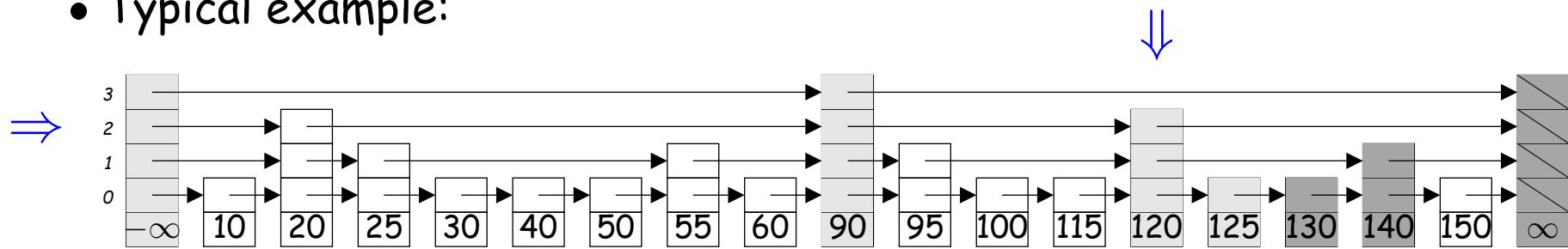
- A *skip list* can be thought of as a kind of n-ary search tree in which we choose to put the keys at "random" heights.
- More often thought of as an ordered list in which one can skip large segments.
- Typical example:



- To search, start at top layer on left, search until next step would overshoot, then go down one layer and repeat.
- In list above, we search for 125 and 127. Gray nodes are looked at; darker gray nodes are overshoots.
- Heights of the nodes were chosen randomly so that there are about $1/2$ as many nodes that are $> k$ high as there are that are k high.
- Makes searches fast *with high probability*.

Probabilistic Balancing: Skip Lists

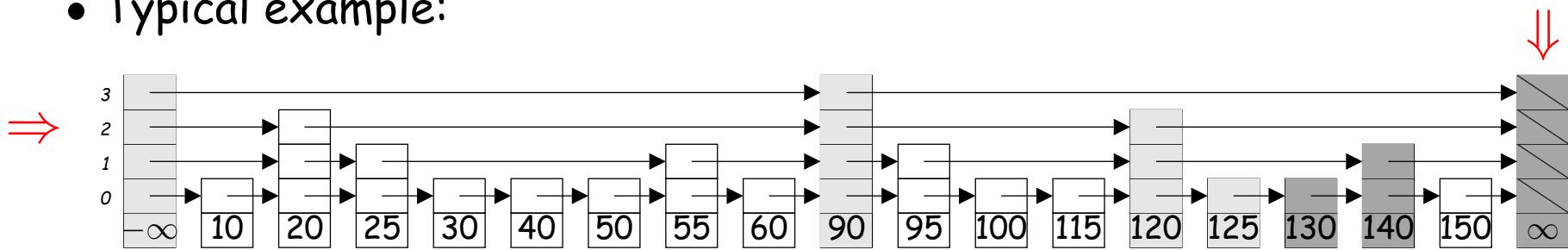
- A *skip list* can be thought of as a kind of n-ary search tree in which we choose to put the keys at "random" heights.
- More often thought of as an ordered list in which one can skip large segments.
- Typical example:



- To search, start at top layer on left, search until next step would overshoot, then go down one layer and repeat.
- In list above, we search for 125 and 127. Gray nodes are looked at; darker gray nodes are overshoots.
- Heights of the nodes were chosen randomly so that there are about $1/2$ as many nodes that are $> k$ high as there are that are k high.
- Makes searches fast *with high probability*.

Probabilistic Balancing: Skip Lists

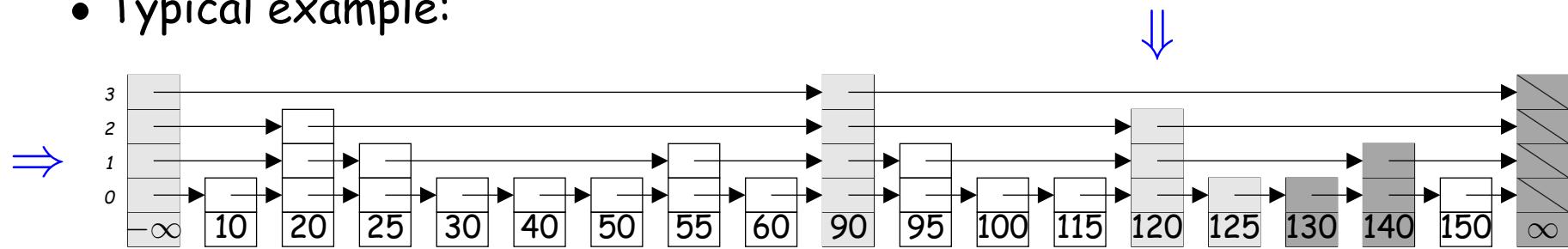
- A *skip list* can be thought of as a kind of n-ary search tree in which we choose to put the keys at "random" heights.
- More often thought of as an ordered list in which one can skip large segments.
- Typical example:



- To search, start at top layer on left, search until next step would overshoot, then go down one layer and repeat.
- In list above, we search for 125 and 127. Gray nodes are looked at; darker gray nodes are overshoots.
- Heights of the nodes were chosen randomly so that there are about $1/2$ as many nodes that are $> k$ high as there are that are k high.
- Makes searches fast *with high probability*.

Probabilistic Balancing: Skip Lists

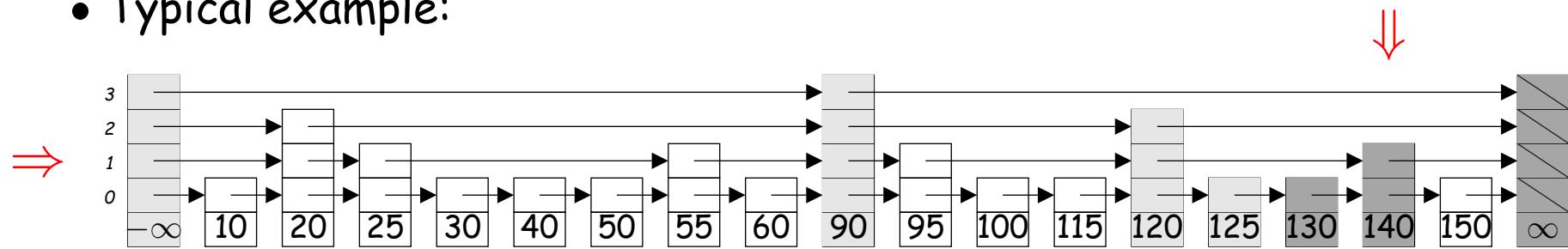
- A *skip list* can be thought of as a kind of n-ary search tree in which we choose to put the keys at "random" heights.
- More often thought of as an ordered list in which one can skip large segments.
- Typical example:



- To search, start at top layer on left, search until next step would overshoot, then go down one layer and repeat.
- In list above, we search for 125 and 127. Gray nodes are looked at; darker gray nodes are overshoots.
- Heights of the nodes were chosen randomly so that there are about $1/2$ as many nodes that are $> k$ high as there are that are k high.
- Makes searches fast *with high probability*.

Probabilistic Balancing: Skip Lists

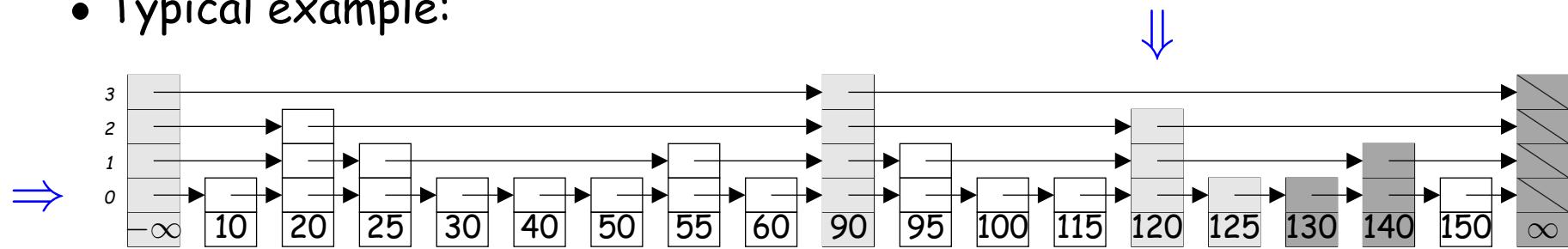
- A *skip list* can be thought of as a kind of n-ary search tree in which we choose to put the keys at "random" heights.
- More often thought of as an ordered list in which one can skip large segments.
- Typical example:



- To search, start at top layer on left, search until next step would overshoot, then go down one layer and repeat.
- In list above, we search for 125 and 127. Gray nodes are looked at; darker gray nodes are overshoots.
- Heights of the nodes were chosen randomly so that there are about $1/2$ as many nodes that are $> k$ high as there are that are k high.
- Makes searches fast *with high probability*.

Probabilistic Balancing: Skip Lists

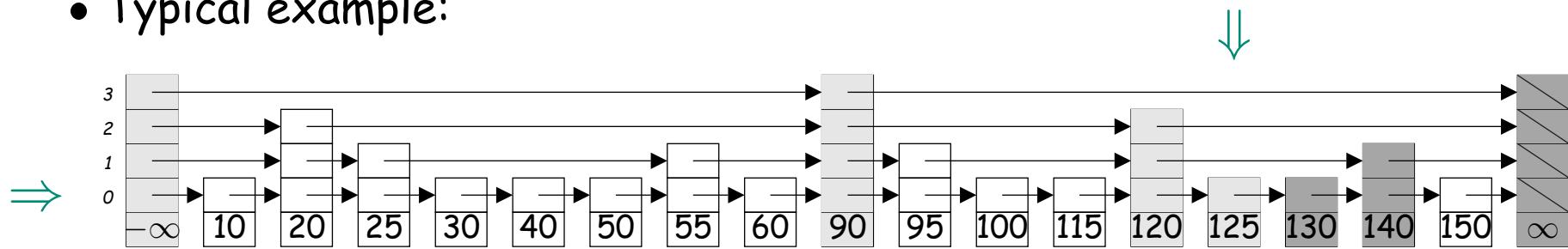
- A *skip list* can be thought of as a kind of n-ary search tree in which we choose to put the keys at "random" heights.
- More often thought of as an ordered list in which one can skip large segments.
- Typical example:



- To search, start at top layer on left, search until next step would overshoot, then go down one layer and repeat.
- In list above, we search for 125 and 127. Gray nodes are looked at; darker gray nodes are overshoots.
- Heights of the nodes were chosen randomly so that there are about $1/2$ as many nodes that are $> k$ high as there are that are k high.
- Makes searches fast *with high probability*.

Probabilistic Balancing: Skip Lists

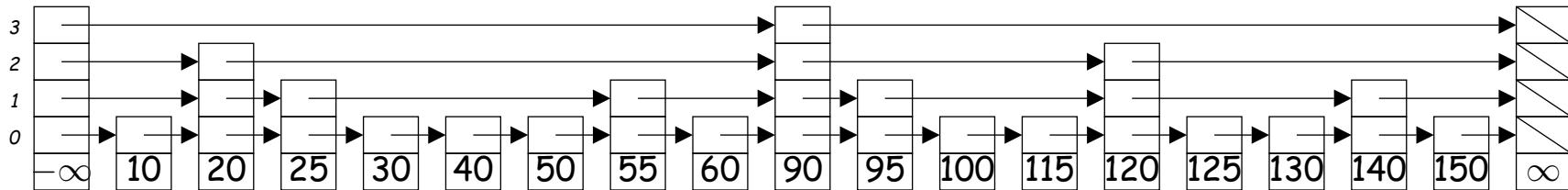
- A *skip list* can be thought of as a kind of n-ary search tree in which we choose to put the keys at "random" heights.
- More often thought of as an ordered list in which one can skip large segments.
- Typical example:



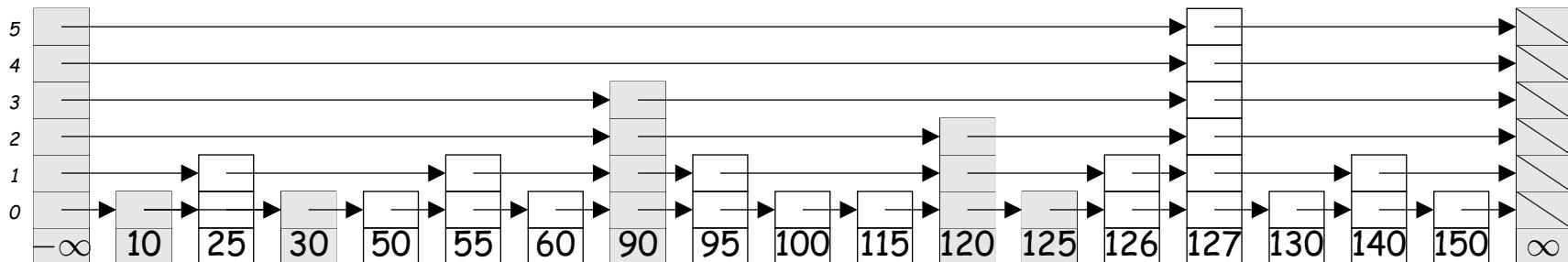
- To search, start at top layer on left, search until next step would overshoot, then go down one layer and repeat.
- In list above, we search for 125 and 127. Gray nodes are looked at; darker gray nodes are overshoots.
- Heights of the nodes were chosen randomly so that there are about $1/2$ as many nodes that are $> k$ high as there are that are k high.
- Makes searches fast *with high probability*.

Example: Adding and deleting

- Starting from initial list:



- In any order, we add 126 and 127 (choosing random heights for them), and remove 20 and 40:

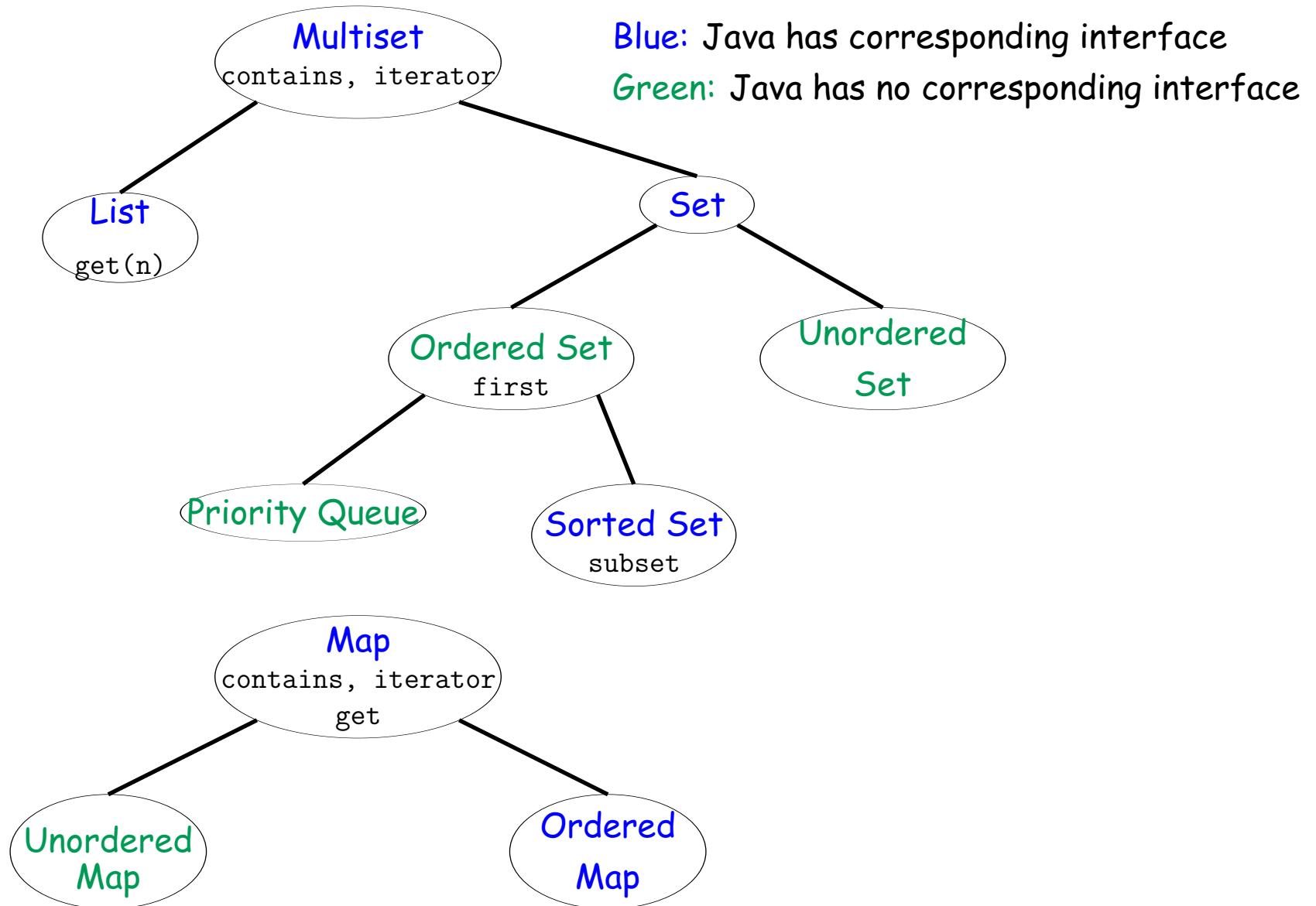


- Shaded nodes here have been modified.

Summary

- Balance in search trees allows us to realize $\Theta(\lg N)$ performance.
- B-trees, red-black trees:
 - Give $\Theta(\lg N)$ performance for searches, insertions, deletions.
 - B-trees good for external storage. Large nodes minimize # of I/O operations
- Tries:
 - Give $\Theta(B)$ performance for searches, insertions, and deletions, where B is length of key being processed.
 - But hard to manage space efficiently.
- *Interesting idea:* scrunched arrays share space.
- Skip lists:
 - Give probable $\Theta(\lg N)$ performance for searches, insertions, deletions
 - Easy to implement.
 - Presented for *interesting ideas:* probabilistic balance, randomized data structures.

Summary of Collection Abstractions



Data Structures that Implement Abstractions

Multiset

- **List**: arrays, linked lists, circular buffers
- **Set**
 - **Ordered Set**
 - * **Priority Queue**: heaps
 - * **Sorted Set**: binary search trees, red-black trees, B-trees, sorted arrays or linked lists
 - **Unordered Set**: hash table

Map

- **Unordered Map**: hash table
- **Ordered Map**: red-black trees, B-trees, sorted arrays or linked lists

Corresponding Classes in Java

Multiset (Collection)

- List: ArrayList, LinkedList, Stack, ArrayBlockingQueue, ArrayDeque
- Set
 - OrderedSet
 - * Priority Queue: PriorityQueue
 - * Sorted Set (SortedSet): TreeSet
 - Unordered Set: HashSet

Map

- Unordered Map: HashMap
- Ordered Map (SortedMap): TreeMap

Lecture #32

Git: A Case Study in System and Data-Structure Design

- Git is a distributed version-control system, apparently the most popular of these currently.
- Conceptually, it stores snapshots (*versions*) of the files and directory structure of a project, keeping track of their relationships, authors, dates, and log messages.
- It is *distributed*, in that there can be many copies of a given repository, each supporting independent development, with machinery to transmit and reconcile versions between repositories.
- Its operation is extremely fast (as these things go).

A Little History

- Developed by Linus Torvalds and others in the Linux community when the developer of their previous, proprietary VCS (Bitkeeper) withdrew the free version.
- The initial implementation effort seems to have taken about 2-3 months, in time for the 2.6.12 Linux kernel release in June, 2005.
- As for the name, according to Wikipedia,

Torvalds has quipped about the name Git, which is British English slang meaning "unpleasant person". Torvalds said: "I'm an egotistical bastard, and I name all my projects after myself. First 'Linux', now 'git'." The man page describes Git as "the stupid content tracker."

- Initially, it was a collection of basic primitives (now called "plumbing") that could be scripted together to provide the desired functionality.
- Then, higher-level commands ("porcelain") was built on top of these to provide a convenient user interface.

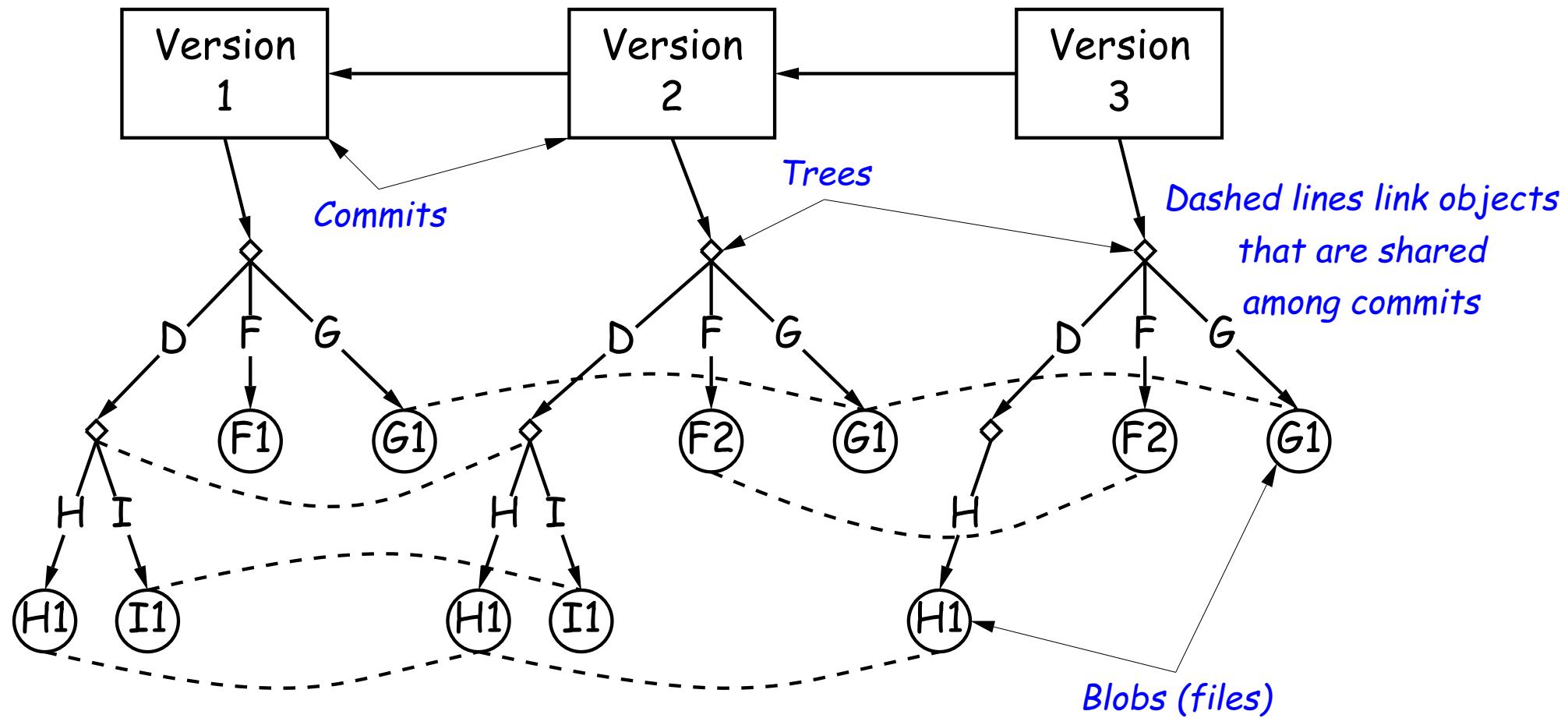
User-Level Conceptual Structure

- The main abstraction is that of a graph of versions or snapshots (called *commits*) of a complete project.
- The graph structure reflects ancestry: which versions came from which.
- Each commit contains
 - A directory tree of files (like a Unix directory).
 - Information about who committed and when.
 - Log message.
 - Pointers to commit (or commits, if there was a merge) from which the commit was derived.

Internal Structures

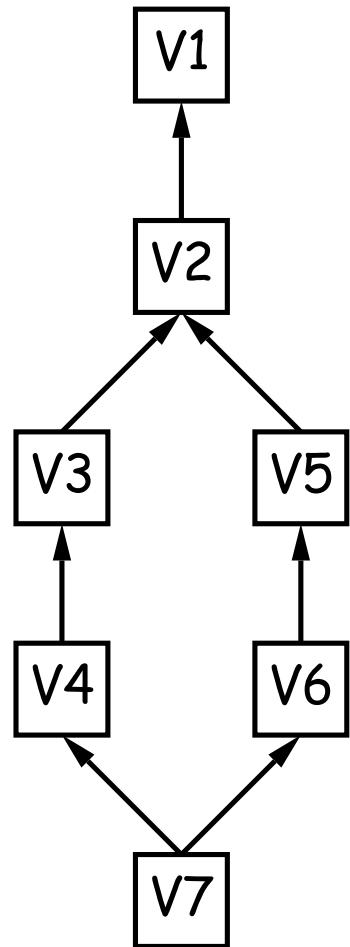
- The main internal components consist of four types of *object*:
 - *Blobs*: basically files of text or bytes.
 - *Trees*: directory structures of files.
 - *Commits*: objects containing references to trees, ancestor commits, and additional information (committer, date, log message).
 - *Tags*: references to commits or other objects, with additional information, intended to identify releases, other important versions, or various useful information. (We won't go into this today).

Commits, Trees, Files

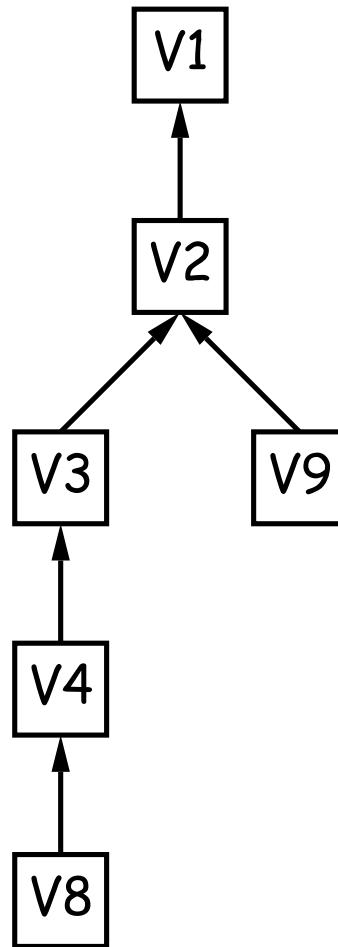


Version Histories in Two Repositories

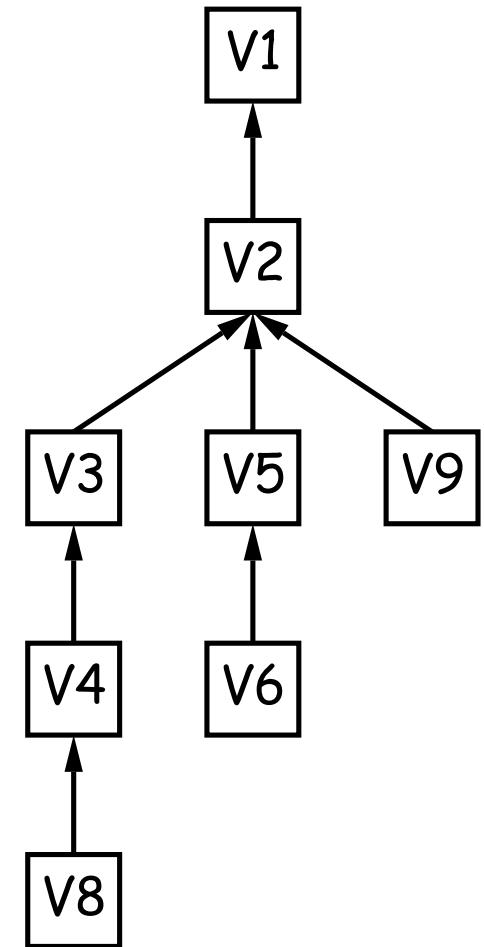
Repository 1



Repository 2



Repository 2
after pushing V6 to it



Major User-Level Features (II)

- Each commit has a name that uniquely identifies it to all repositories.
- Repositories can transmit collections of versions to each other.
- Transmitting a commit from repository A to repository B requires only the transmission of those objects (files or directory trees) that B does not yet have (allowing speedy updating of repositories).
- Repositories maintain named **branches**, which are simply identifiers of particular commits that are updated to keep track of the most recent commits in various lines of development.
- Likewise, **tags** are essentially named pointers to particular commits. Differ from branches in that they are not usually changed.

Internals

- Each Git repository is contained in a directory.
- A repository may either be *bare* (just a collection of objects and metadata), or it may be included as part of a working directory (where it is named `.git`).
- The data of the repository is stored in various *objects* represented as files corresponding to data files (or other “leaf” content), trees, and commits.
- To save space, data in these objects is *compressed*.
- Git can *garbage-collect* the objects from time to time to save additional space and can *pack* bunches of objects into compressed files.

The Pointer Problem

- Objects in Git are represented as files. How should we represent pointers between them (e.g., from trees to files, or from commits to trees and other commits)?
- We want to be able to *transmit* objects from one repository to another with different contents. How do you transmit the pointers?
 - ... Since the address of an object on one machine is very unlikely to have anything to do with its address on another.
- We only want to transfer those objects that are missing in the target repository. How do we know which those are?
- We could use a counter in each repository to give each object there a unique file name. But how can that work consistently for two independent repositories?

Content-Addressable File System

- Could use some way of naming objects that is universal.
- We use the names, then, as pointers.
- Solves the “Which objects don’t you have?” problem in an obvious way.
- Conceptually, what is invariant about a version of an object, regardless of repository, is its **contents**.
- But we can’t use the contents as the name for obvious reasons (like trying to use a country as its own map.)
- **Idea:** Use a **hash of the contents** as the address.
- **Problem:** That doesn’t work!

Content-Addressable File System

- Could use some way of naming objects that is universal.
- We use the names, then, as pointers.
- Solves the “Which objects don’t you have?” problem in an obvious way.
- Conceptually, what is invariant about a version of an object, regardless of repository, is its **contents**.
- But we can’t use the contents as the name for obvious reasons (like trying to use a country as its own map.)
- **Idea:** Use a **hash of the contents** as the address.
- **Problem:** That doesn’t work!
- **Brilliant Idea:** Use it anyway!!

How A Broken Idea Can Work

- The idea is to use a hash function that is so unlikely to have a collision that we can ignore that possibility.
- *Cryptographic Hash Functions* have relevant properties.
- Such a function, f , is designed to withstand cryptoanalytic attacks. In particular, it should have
 - *Pre-image resistance*: given $h = f(m)$, it should be computationally infeasible to find such a message m , given h .
 - *Second pre-image resistance*: given message m_1 , it should be infeasible to find $m_2 \neq m_1$ such that $f(m_1) = f(m_2)$.
 - *Collision resistance*: it should be difficult to find *any* two messages $m_1 \neq m_2$ such that $f(m_1) = f(m_2)$.
- With these properties, the scheme of using hashes of contents as names is extremely unlikely to fail, even when the system is used maliciously.

SHA1

- Git uses *SHA1* (Secure Hash Function 1).
- Can play around with this using the `hashlib` module in Python3.
- All object names in Git are therefore 160-bit hash codes of contents, in hex.
- E.g. a recent commit in the Spring 2022 shared CS61B repository could be fetched (if needed) with

```
git checkout 72c136ba61fb468bf4c8ac906cd57a2c7fa25025
```

Low-Level Blob Management

- You can find out the hashcode that Git uses for the blob containing file something.java with the command

```
git hash-object something.java
```

- And if this tells you that the file would have hash code

192a0ca0d159f1550b0b5e102f7e06867cc44782

and you actually git add this file, its compressed contents will be stored (as a blob) in the file

.git/objects/19/2a0ca0d159f1550b0b5e102f7e06867cc44782

[Why do you suppose the implementors of Git chose this name?]

- You can look at the (uncompressed) contents of the blob with

```
git cat-file -p 192a0ca0d159f1550b0b5e102f7e06867cc44782
```

CS61B Lecture #33

Today's Readings: Graph Structures: *DSIJ*, Chapter 12

Why Graphs?

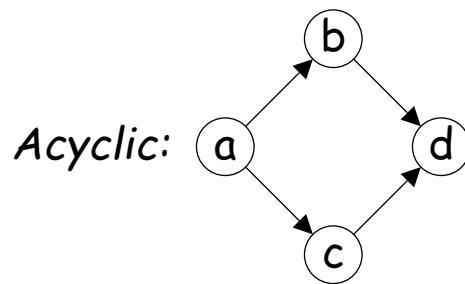
- For expressing non-hierarchically related items
- Examples:
 - Networks: pipelines, roads, assignment problems
 - Representing processes: flow charts, Markov models
 - Representing partial orderings: PERT charts, makefiles
 - As we've seen, in representing connected structures as used in Git.

Some Terminology

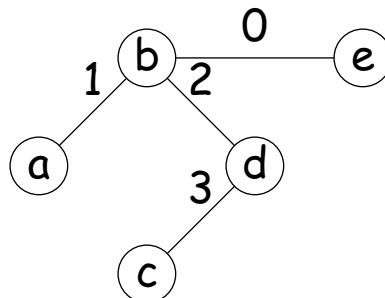
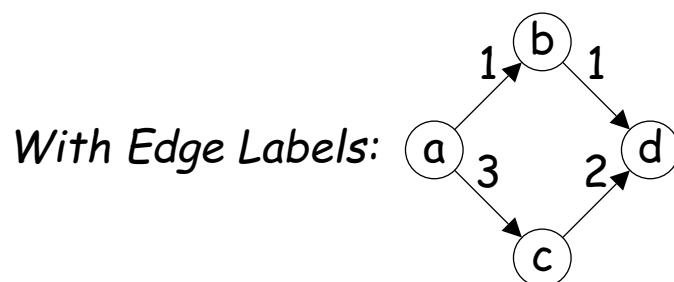
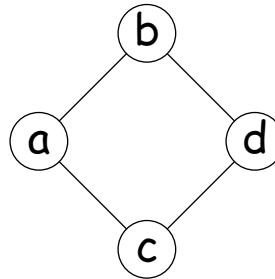
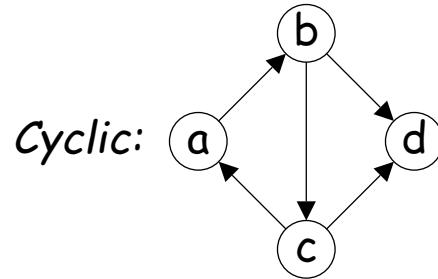
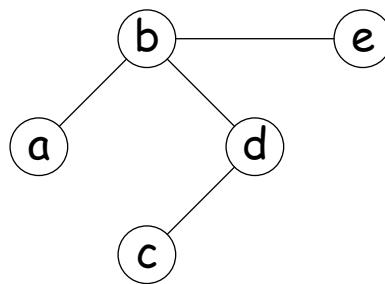
- A *graph* consists of
 - A set of *nodes* (aka *vertices*)
 - A set of *edges*: pairs of nodes.
 - Nodes with an edge between are *adjacent*.
 - Depending on problem, nodes or edges may have *labels* (or *weights*)
- Typically call the node set $V = \{v_0, \dots\}$, and the edge set E .
- If the edges have an order (first, second), they are *directed edges*, and we have a *directed graph (digraph)*; otherwise an *undirected graph*.
- Edges are *incident* to their nodes.
- Directed edges *exit* one node and *enter* the next.
- A *cycle* is a *path*—a sequence of edges—with repeated edges leading from a node back to itself (following arrows if directed).
- A graph is *cyclic* if it has a cycle; else *acyclic*. Abbreviation: Directed Acyclic Graph—**DAG**.

Some Pictures

Directed

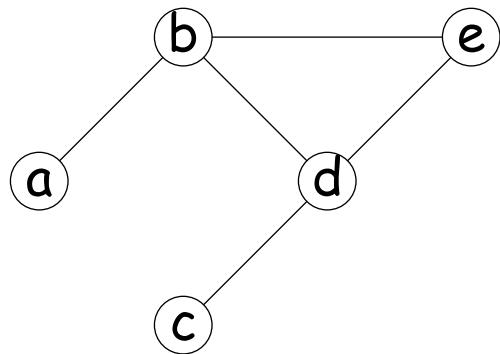


Undirected

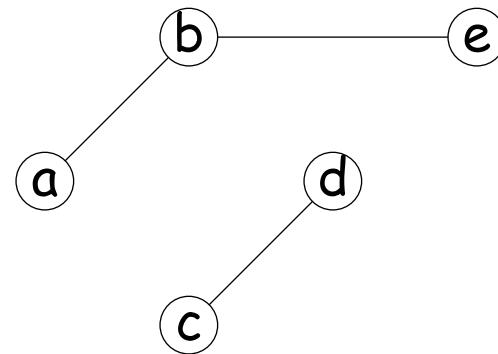


Connectivity (Undirected)

- A *path* is variously defined as a sequence of vertices v_0, v_1, \dots, v_n where there is an edge from each v_i to v_{i+1} , or as the sequence of edges between them: $(v_0, v_1), (v_1, v_2), \dots$.
- An undirected graph is *connected* if there is a *path* between every pair of nodes in the graph:



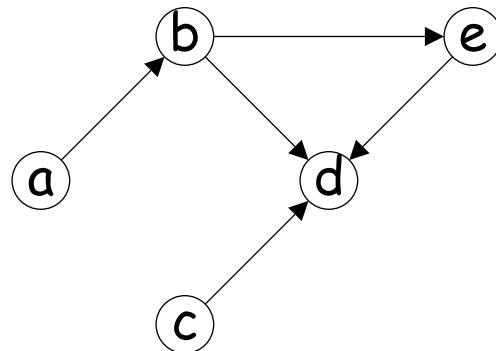
Connected



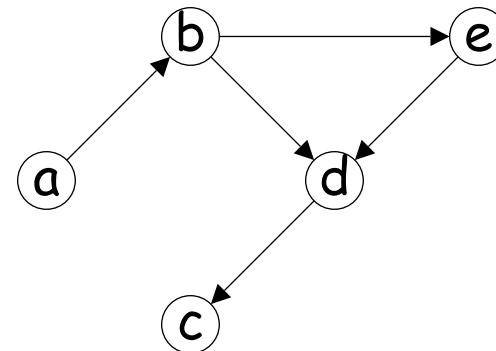
Unconnected

Connectivity (Directed)

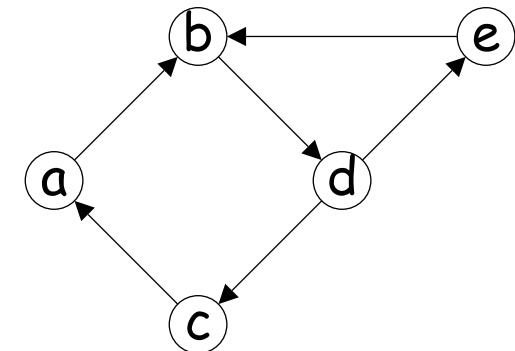
- For directed graphs, it's more complicated:
 - *Weakly connected*: connected if direction is removed.
 - *Unilaterally connected* (or *semiconnected*): there is a path in one direction or the other between each node pair;
 - *Strongly connected*: there are paths in both directions for each node pair.



Weakly Connected



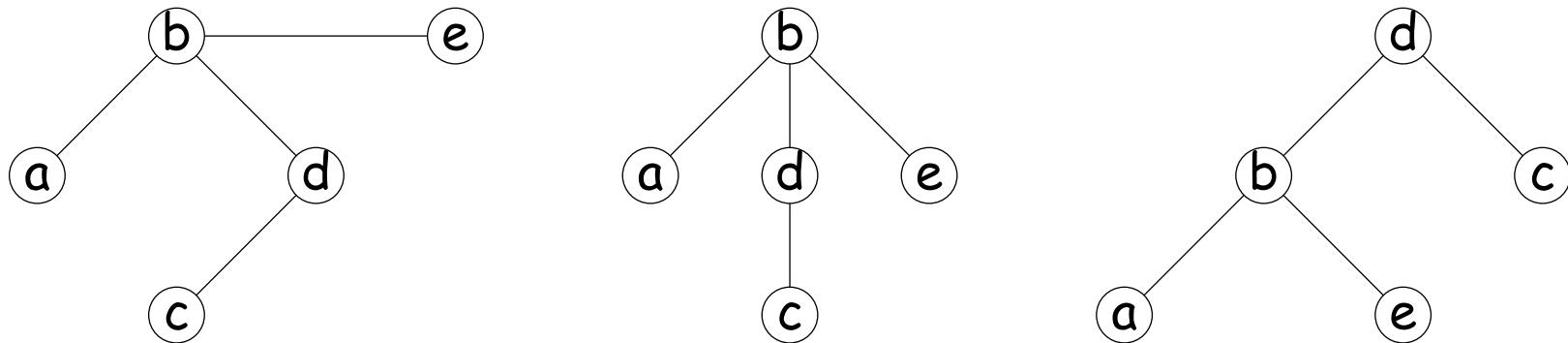
Unilaterally Connected



Strongly Connected

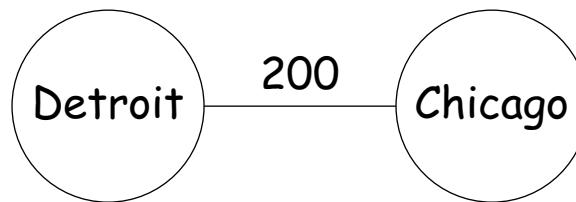
Trees and Graphs

- A DAG is a (rooted) tree iff connected, and every node but the root has exactly one parent.
- A connected, acyclic, undirected graph is also called a *free tree*.
Free: we're free to pick the root; e.g., all the following are the same graph:

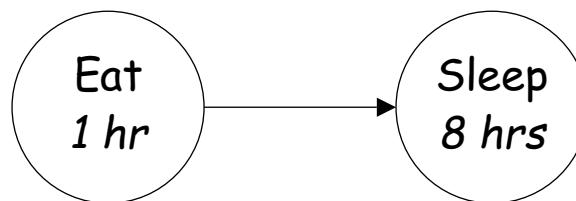


Examples of Use

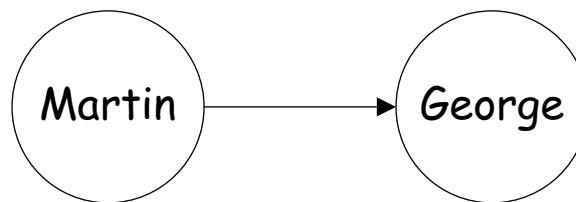
- Edge = Connecting road, with length.



- Edge = Must be completed before; Node label = time to complete.

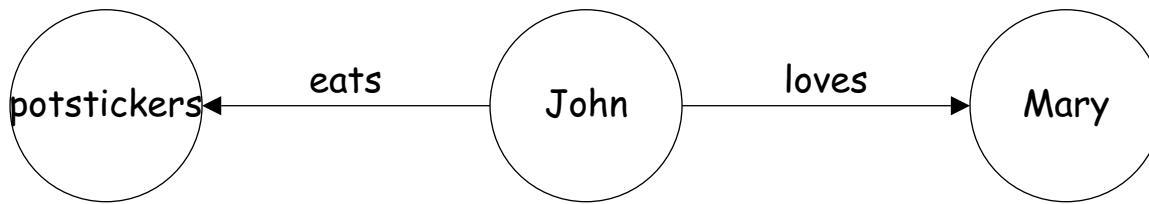


- Edge = Begat

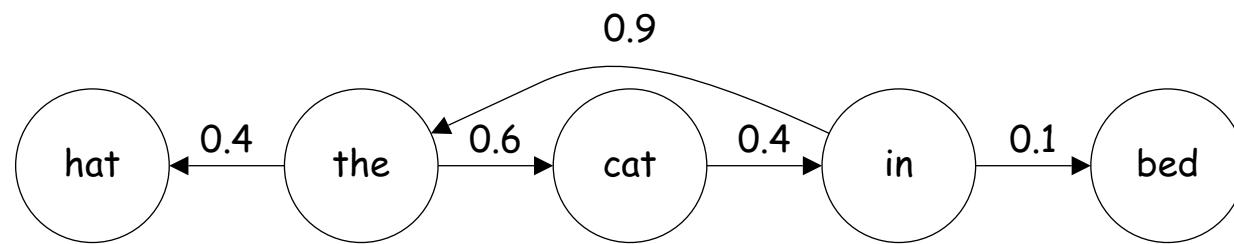


More Examples

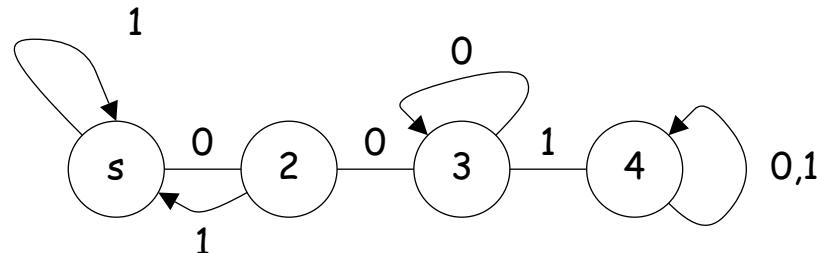
- Edge = some relationship



- Edge = next state might be (with probability)

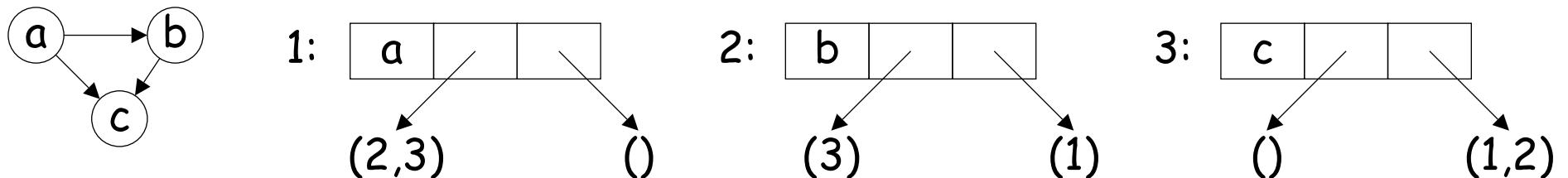


- Edge = next state in state machine, label is triggering input. (Start at s. Being in state 4 means “there is a substring ‘001’ somewhere in the input”.)



Representation

- Often useful to number the nodes, and use the numbers in edges.
- *Edge list representation*: each node contains some kind of list (e.g., linked list or array) of its successors (and possibly predecessors).



- *Edge sets*: Collection of all edges. For graph above:

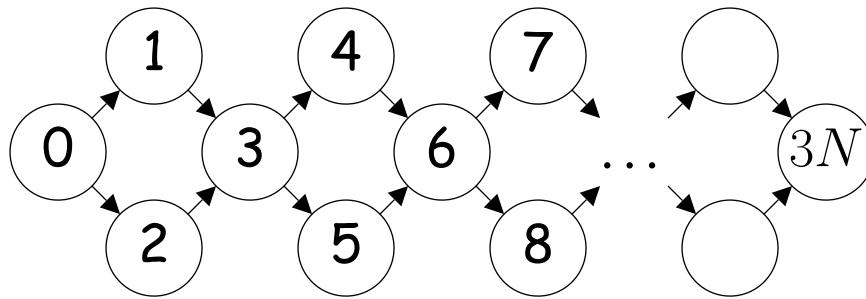
$$\{(1, 2), (1, 3), (2, 3)\}$$

- *Adjacency matrix*: Represent connection with matrix entry:

$$\begin{array}{c} & \begin{matrix} 1 & 2 & 3 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \end{matrix} & \left[\begin{matrix} 0 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{matrix} \right] \end{array}$$

Traversing a Graph

- Many algorithms on graphs depend on traversing all or some nodes.
- Can't quite use recursion because of cycles.
- Even in acyclic graphs, can get combinatorial explosions:



Treat 0 as the root and do recursive traversal down the two edges out of each node: $\Theta(2^N)$ operations!

- So typically try to visit each node constant # of times (e.g., once).

Recursive Depth-First Traversal of a Graph

- Can fix looping and combinatorial problems using the “bread-crumb” method used in earlier lectures for a maze.
- That is, *mark* nodes as we traverse them and don’t traverse previously marked nodes.
- Makes sense to talk about *preorder* and *postorder*, as for trees.

```
void preorderTraverse(Graph G, Node v)
{
    if (v is unmarked) {
        mark(v);
        visit v;
        for (Edge(v, w) ∈ G)
            traverse(G, w);
    }
}

void postorderTraverse(Graph G, Node v)
{
    if (v is unmarked) {
        mark(v);
        for (Edge(v, w) ∈ G)
            traverse(G, w);
        visit v;
    }
}
```

Recursive Depth-First Traversal of a Graph (II)

- We are often interested in traversing *all* nodes of a graph, not just those reachable from one node.
- So we can repeat the procedure as long as there are unmarked nodes.

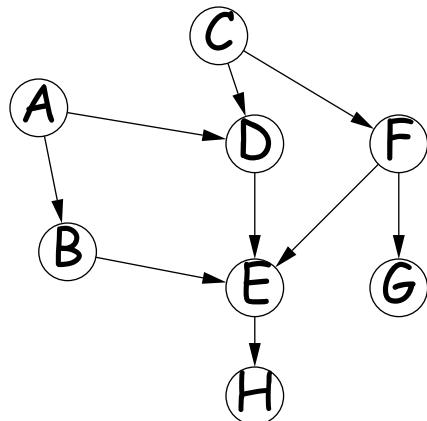
```
void preorderTraverse(Graph G) {  
    clear all marks;  
    for (v ∈ nodes of G) {  
        preorderTraverse(G, v);  
    }  
}  
  
void postorderTraverse(Graph G) {  
    clear all marks;  
    for (v ∈ nodes of G) {  
        postorderTraverse(G, v);  
    }  
}
```

Topological Sorting

Problem: Given a DAG, find a linear order of nodes consistent with the edges.

- That is, order the nodes v_0, v_1, \dots such that v_k is never reachable from $v_{k'}$ if $k' > k$.
- Gmake does this. Also PERT charts.

Graph (two views)

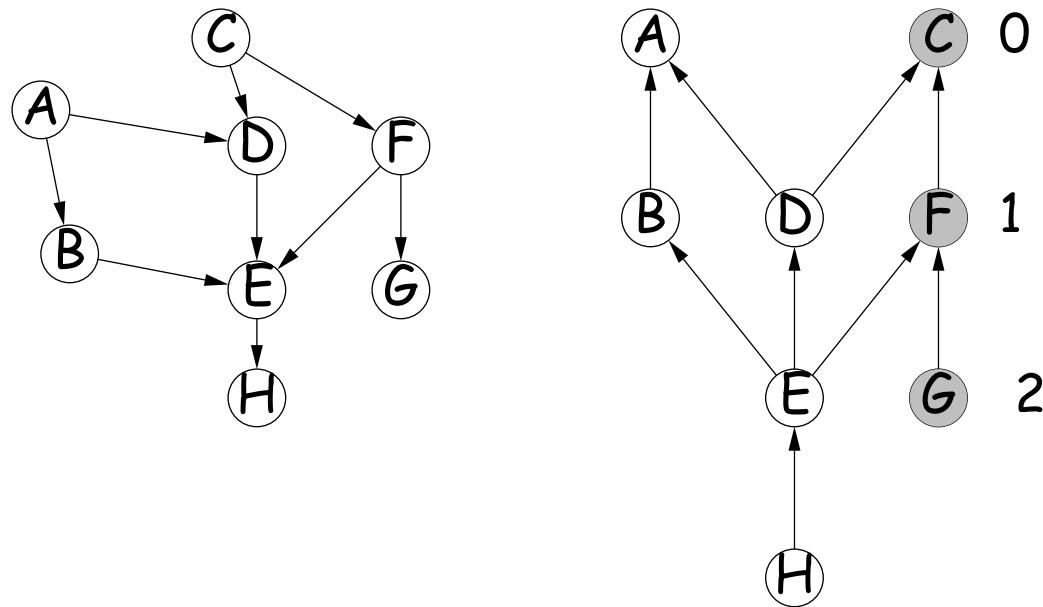


Possible Orderings

A	C	C
C	A	F
B	F	G
D	D	A
F	B	B
E	G	D
G	E	E
H	H	H

Sorting and Depth First Search

- **Observation:** Suppose we *reverse the links* on our graph.
- If we do a recursive DFS on the reverse graph, starting from node H, for example, we will find all nodes that must come *before* H.
- When the search reaches a node in the reversed graph and there are no successors, we know that it is safe to put that node first.
- In general, a *postorder* traversal of the *reversed* graph visits nodes only after all predecessors have been visited.



Numbers show post-order traversal order starting from G: everything that must come before G.

General Graph Traversal Algorithm

```
COLLECTION_OF_VERTICES fringe;  
  
fringe = INITIAL_COLLECTION;  
while (!fringe.isEmpty()) {  
    Vertex v = fringe.REMOVE_HIGHEST_PRIORITY_ITEM();  
  
    if (!MARKED(v)) {  
        MARK(v);  
        VISIT(v);  
        For each edge(v,w) {  
            if (NEEDS_PROCESSING(w))  
                Add w to fringe;  
        }  
    }  
}
```

Replace *COLLECTION_OF_VERTICES*, *INITIAL_COLLECTION*, etc. with various types, expressions, or methods to different graph algorithms.

Example: Depth-First Traversal

Problem: Visit every node reachable from v once, visiting nodes further from start first.

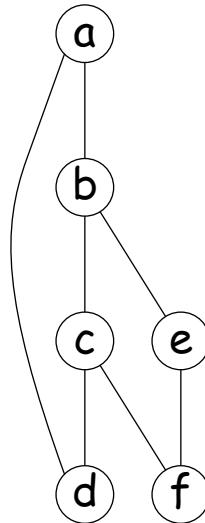
```
// Red sections are specializations of general algorithm
Stack<Vertex> fringe;

fringe = stack containing {v};
while (!fringe.isEmpty()) {
    Vertex v = fringe.pop();

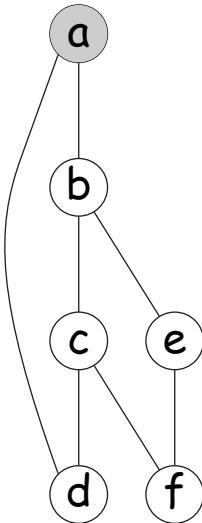
    if (!marked(v)) {
        mark(v);
        VISIT(v);
        For each edge(v,w) {
            if (!marked(w))
                fringe.push(w);
        }
    }
}
```

Depth-First Traversal Illustrated

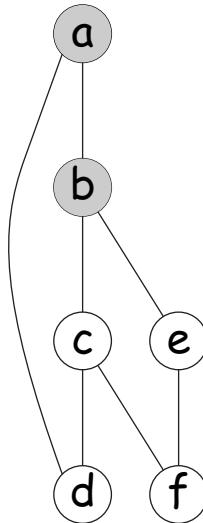
Marked: z



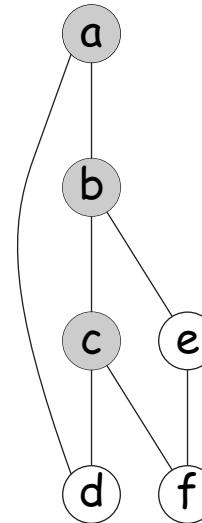
Fringe: [a]



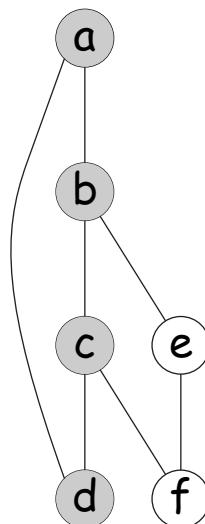
Fringe: [b, d]



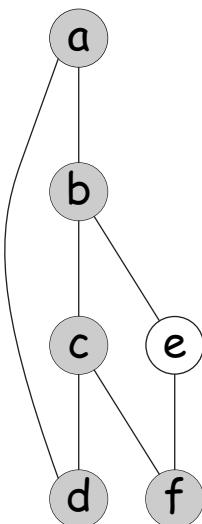
Fringe: [c, e, d]



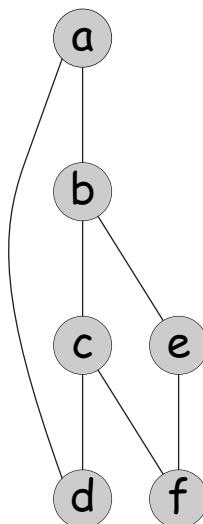
Fringe: [d, f, e, d]



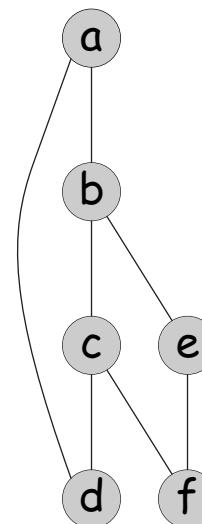
[f, e, d]



[e, e, d]



[e, d]



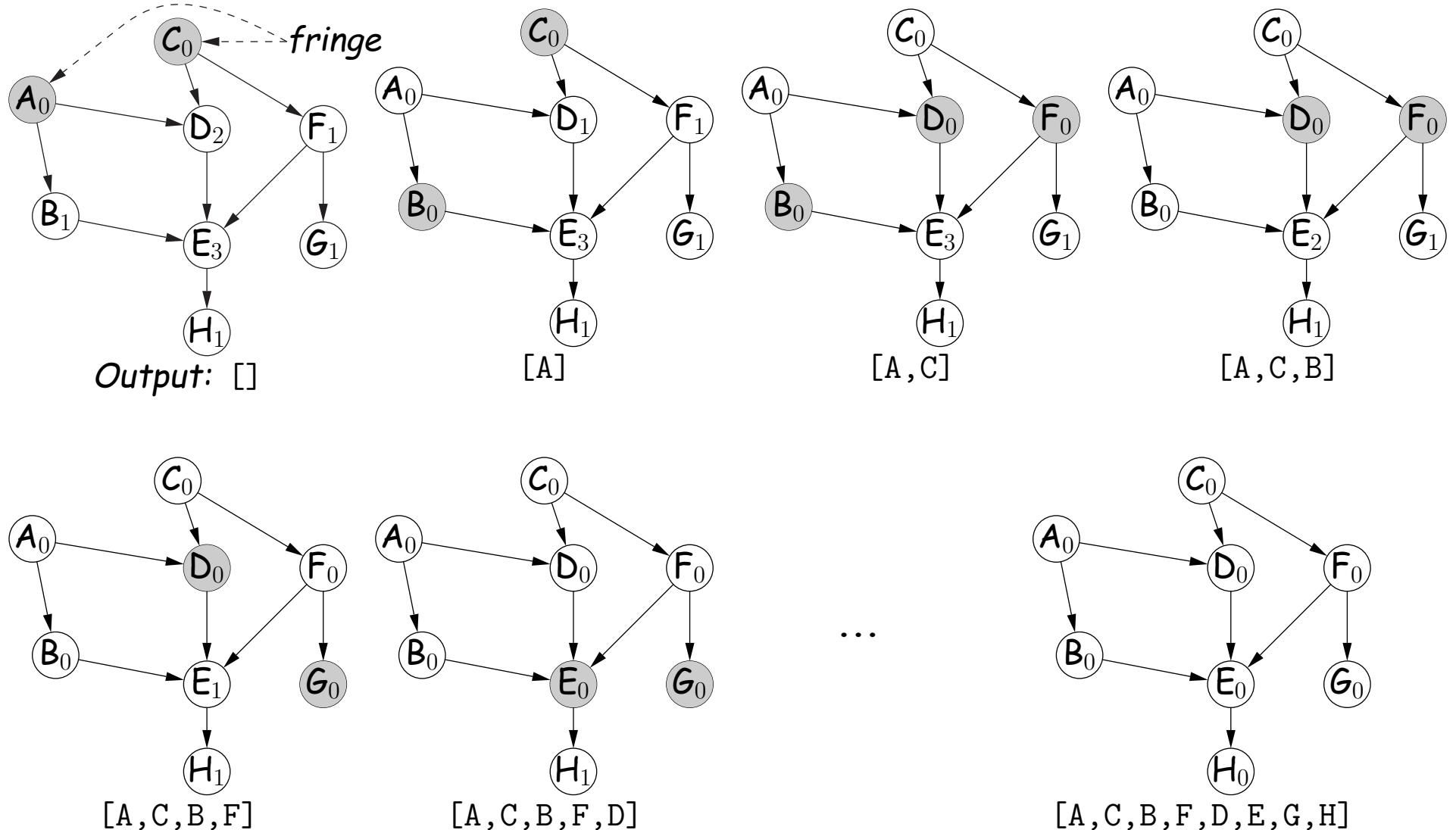
[]

Topological Sort Revisited

- Another approach to topological sorting uses the general traversal algorithm.
- We keep a sequence of output nodes, initialized to empty.
- For each node, maintain a count of the number of immediate predecessor nodes (the number of nodes with edges going to that node) that have not yet been output.
- At each step, output a node with no remaining predecessors, and update the predecessor counts.

```
Set<Vertex> fringe;
Map<Vertex, Integer> predCount;
fringe = empty set of vertices;
predCount = map from each vertex to number of predecessors;
while (!fringe.isEmpty()) {
    Vertex v = fringe.removeSomeElement(); /* Not a real operation! */
    Add v to the output;
    for each edge(v,w) {
        predCount.put(w, predCount.get(w) - 1);
        if (predCount.get(w) == 0) fringe.add(w);
    }
}
```

Alternative Topological Sort in Action



Shortest Paths: Dijkstra's Algorithm

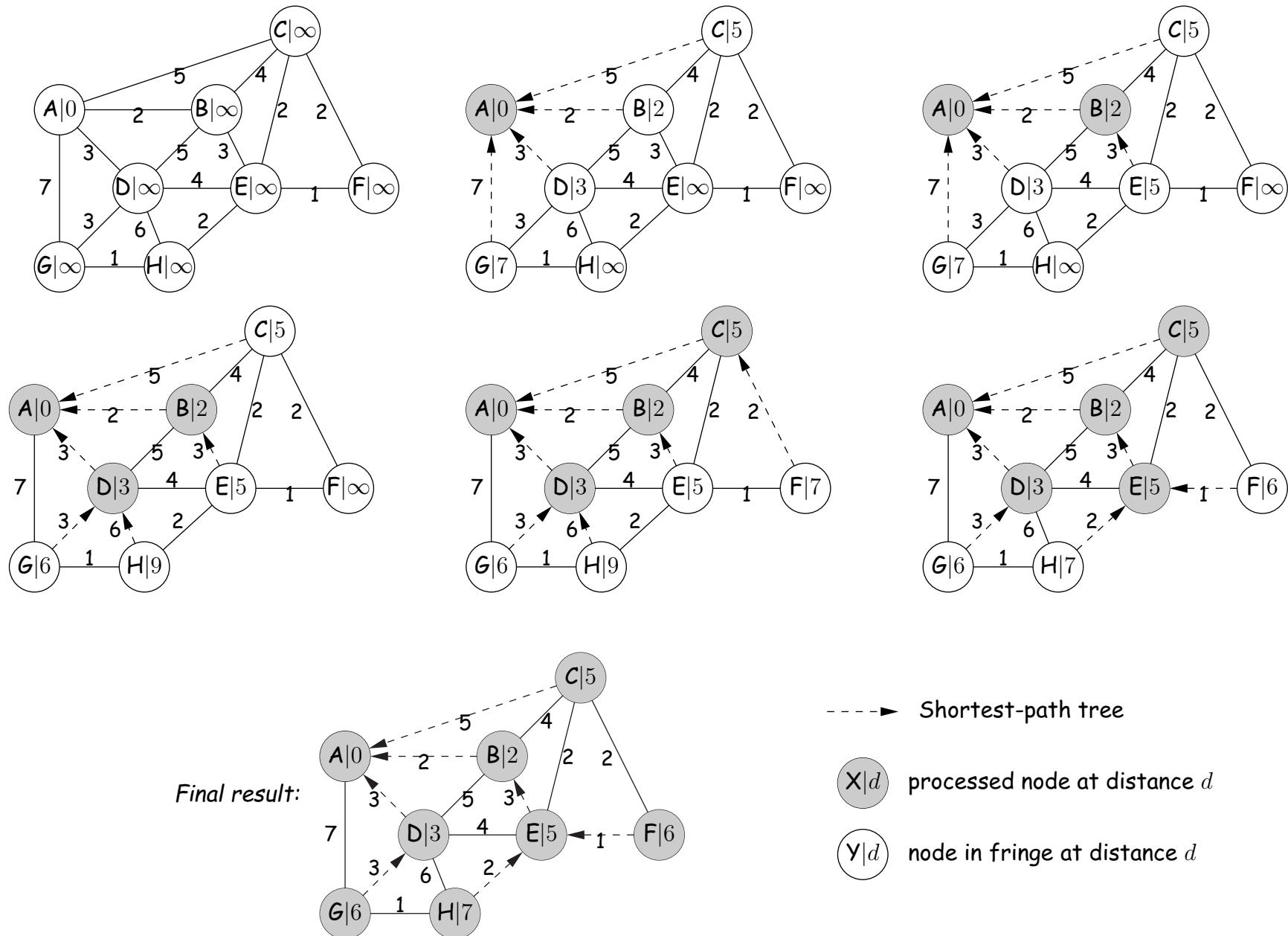
Problem: Given a graph (directed or undirected) with non-negative edge weights, compute shortest paths from given source node, s , to all nodes.

- “Shortest” = sum of weights along path is smallest.
- For each node, keep estimated distance from s , ...
- ... and of preceding node in shortest path from s .

```
PriorityQueue<Vertex> fringe;
For each node v { v.dist() = infinity; v.back() = null; }
s.dist() = 0;
fringe = priority queue ordered by smallest .dist();
add all vertices to fringe;
while (!fringe.isEmpty()) {
    Vertex v = fringe.removeFirst();

    For each edge(v,w) {
        if (v.dist() + weight(v,w) < w.dist())
            { w.dist() = v.dist() + weight(v,w); w.back() = v; }
    }
}
```

Example



CS61B Lecture #34

- Today: A* search, Minimum spanning trees, union-find.
-

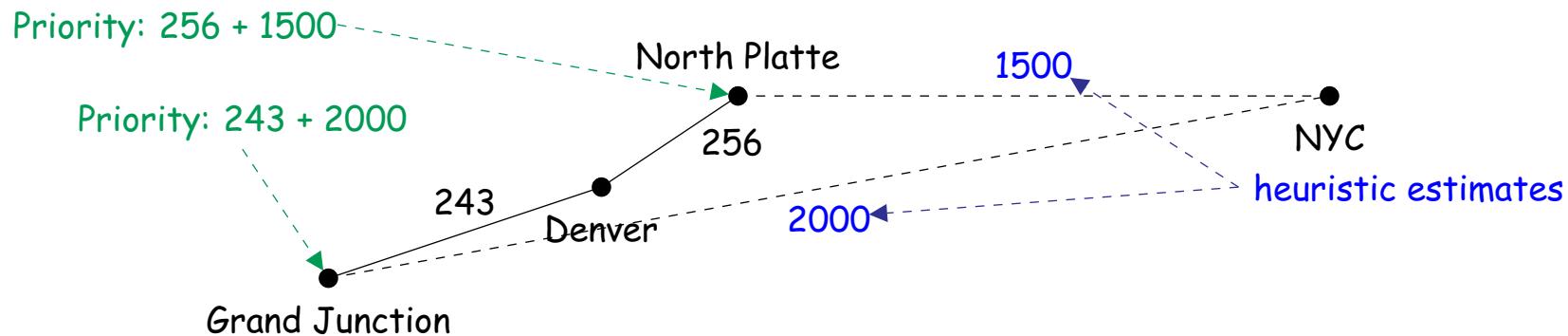
Point-to-Point Shortest Path

- Dijkstra's algorithm gives you shortest paths from a particular given vertex to all others in a graph.
- But suppose you're only interested in getting to a particular vertex?
- Because the algorithm finds paths in order of length, you *could* simply run it and stop when you get to the vertex you want.
- But, this can be really wasteful.
- For example, to travel by road from Denver to a destination on lower Fifth Avenue in New York City is about 1750 miles (says Google).
- But traveling from Denver to Berkeley is about 1250 miles.
- So, we'd explore much of California, Nevada, Arizona, etc. before we found our destination, even though these are all in the wrong direction!
- Situation even worse when graph is infinite, generated on the fly.

A* Search

- We're looking for a path from vertex Denver to the desired NYC vertex.
- Suppose that we had a fast and simple *heuristic estimate*, $h(V)$, of the length of a path from any vertex V to NYC.
- And suppose that instead of visiting vertices in the fringe in order of their shortest known path to Denver, we order by the sum of that distance plus this heuristic estimate of the remaining distance to NYC: $d(\text{Denver}, V) + h(V)$.
- In other words, we look at places that are reachable from places where we already know the shortest path to Denver and choose those that look like they will result in the shortest trip to NYC, guessing at the remaining distance.
- If the estimate is good, then we don't look at, say, Grand Junction (250 miles west by road), because it's in the wrong direction.
- The resulting algorithm is *A* search*.
- But for it to work, we must be careful about the heuristic.

Illustration



- If the solid lines indicate distances from Denver that we have determined so far,
- Dijkstra's algorithm would have us look at Grand Junction next (smallest distance from Denver).
- But if we add in the heuristic remaining distance to NYC (our goal), we choose North Platte instead.

Admissible Heuristics for A* Search

- If our heuristic estimate for the distance to NYC is too high (i.e., larger than the actual path by road), then we may get to NYC without ever examining points along the shortest route.
- For example, if our heuristic decided that the midwest was literally the middle of nowhere, and $h(C) = 2000$ for C any city in Michigan or Indiana, we'd only find a path that detoured south through Kentucky.
- So to be **admissible**, $h(C)$ must never overestimate $d(C, \text{NYC})$, the minimum path distance from C to NYC.
- On the other hand, $h(C) = 0$ will work (what is the result?), but yield a non-optimal algorithm.

Admissible Heuristics for A* Search

- If our heuristic estimate for the distance to NYC is too high (i.e., larger than the actual path by road), then we may get to NYC without ever examining points along the shortest route.
- For example, if our heuristic decided that the midwest was literally the middle of nowhere, and $h(C) = 2000$ for C any city in Michigan or Indiana, we'd only find a path that detoured south through Kentucky.
- So to be **admissible**, $h(C)$ must never overestimate $d(C, \text{NYC})$, the minimum path distance from C to NYC.
- On the other hand, $h(C) = 0$ will work (what is the result?), but yield a non-optimal algorithm. This is just Dijkstra's algorithm.

Consistency

- Suppose that we estimate $h(\text{Chicago}) = 700$, and $h(\text{Springfield, IL}) = 200$, where $d(\text{Chicago}, \text{Springfield}) = 200$.
- So by driving 200 miles to Springfield, we guess that we are suddenly 500 miles closer to NYC.
- This is admissible, since both estimates are low, but it will mess up our algorithm.
- Specifically, will require that we put processed nodes back into the fringe, in case our estimate was wrong.
- So (in this course, anyway) we also require *consistent heuristics*: $h(A) \leq h(B) + d(A, B)$, as for the triangle inequality.
- All consistent heuristics are admissible (why?).
- So distance “as the crow flies” is a reasonable $h(\cdot)$ in the trip-planning applications.
- Demo of A^* search (and others) is in [cs61b-software](#) and on the instructional machines as [graph-demo](#).

Summary of Shortest Paths

- Dijkstra's algorithm finds a *shortest-path tree* computing giving (backwards) shortest paths in a weighted graph from a given starting node to all other nodes.
- Time required =
 - Time to remove V nodes from priority queue +
 - Time to update all neighbors of each of these nodes and add or reorder them in queue ($E \lg E$)
 - $\in \Theta(V \lg V + E \lg V) = \Theta((V + E) \lg V)$
- A* search searches for a shortest path to a *particular* target node.
- Same as Dijkstra's algorithm, except:
 - Stop when we take target from queue.
 - Order queue by estimated distance to start + heuristic guess of remaining distance ($h(v) = d(v, \text{target})$)
 - Heuristic must not overestimate distance and obey triangle inequality ($d(a, b) + d(b, c) \geq d(a, c)$).

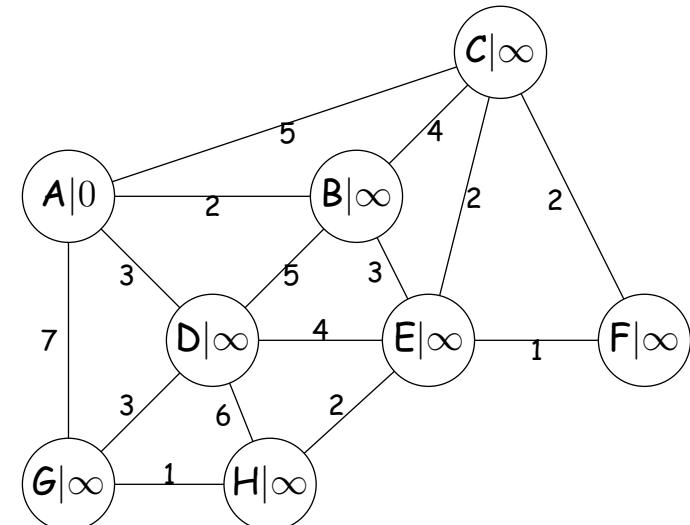
Minimum Spanning Trees

- **Problem:** Given a set of places and distances between them (assume always positive), find a set of connecting roads of minimum total length that allows travel between any two.
- The routes you get will not necessarily be shortest paths.
- Easy to see that such a set of connecting roads and places must form a tree, because removing one road in a cycle still allows all to be reached.

Minimum Spanning Trees by Prim's Algorithm

- Idea is to grow a tree starting from an arbitrary node.
- At each step, add the shortest edge connecting some node already in the tree to one that isn't yet.
- Why must this work?

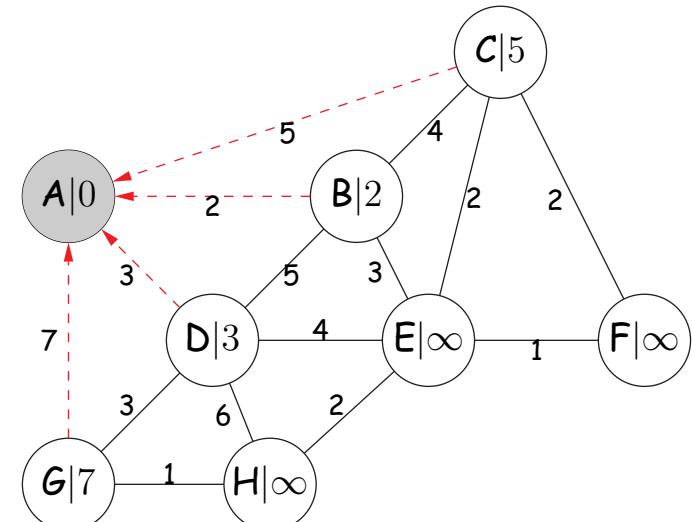
```
PriorityQueue fringe;  
For each node v { v.dist() = infinity; v.parent() = null; }  
Choose an arbitrary starting node, s;  
s.dist() = 0;  
fringe = priority queue ordered by smallest .dist();  
add all vertices to fringe;  
while (!fringe.isEmpty()) {  
    Vertex v = fringe.removeFirst();  
  
    For each edge(v,w) {  
        if (w ∈ fringe && weight(v,w) < w.dist())  
            { w.dist() = weight(v, w); w.parent() = v; }  
    }  
}
```



Minimum Spanning Trees by Prim's Algorithm

- Idea is to grow a tree starting from an arbitrary node.
- At each step, add the shortest edge connecting some node already in the tree to one that isn't yet.
- Why must this work?

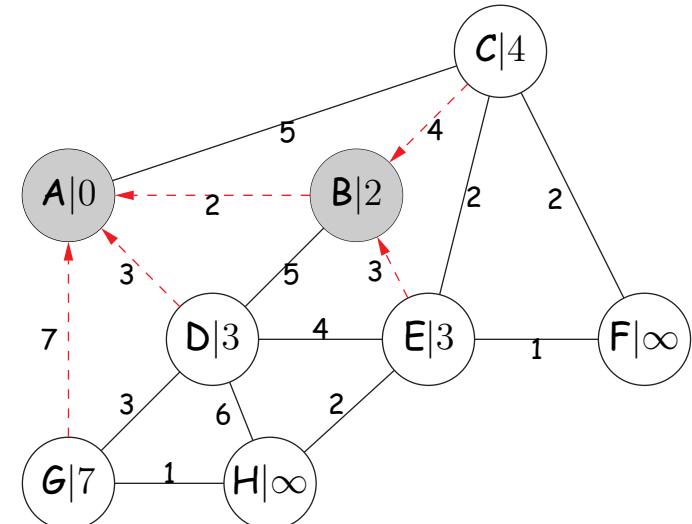
```
PriorityQueue fringe;  
For each node v { v.dist() = infinity; v.parent() = null; }  
Choose an arbitrary starting node, s;  
s.dist() = 0;  
fringe = priority queue ordered by smallest .dist();  
add all vertices to fringe;  
while (!fringe.isEmpty()) {  
    Vertex v = fringe.removeFirst();  
  
    For each edge(v,w) {  
        if (w ∈ fringe && weight(v,w) < w.dist())  
            { w.dist() = weight(v, w); w.parent() = v; }  
    }  
}
```



Minimum Spanning Trees by Prim's Algorithm

- Idea is to grow a tree starting from an arbitrary node.
- At each step, add the shortest edge connecting some node already in the tree to one that isn't yet.
- Why must this work?

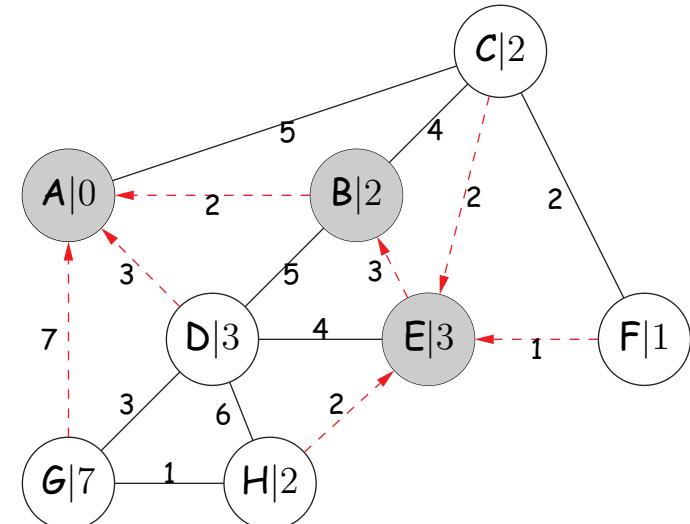
```
PriorityQueue fringe;  
For each node v { v.dist() = infinity; v.parent() = null; }  
Choose an arbitrary starting node, s;  
s.dist() = 0;  
fringe = priority queue ordered by smallest .dist();  
add all vertices to fringe;  
while (!fringe.isEmpty()) {  
    Vertex v = fringe.removeFirst();  
  
    For each edge(v,w) {  
        if (w ∈ fringe && weight(v,w) < w.dist())  
            { w.dist() = weight(v, w); w.parent() = v; }  
    }  
}
```



Minimum Spanning Trees by Prim's Algorithm

- Idea is to grow a tree starting from an arbitrary node.
- At each step, add the shortest edge connecting some node already in the tree to one that isn't yet.
- Why must this work?

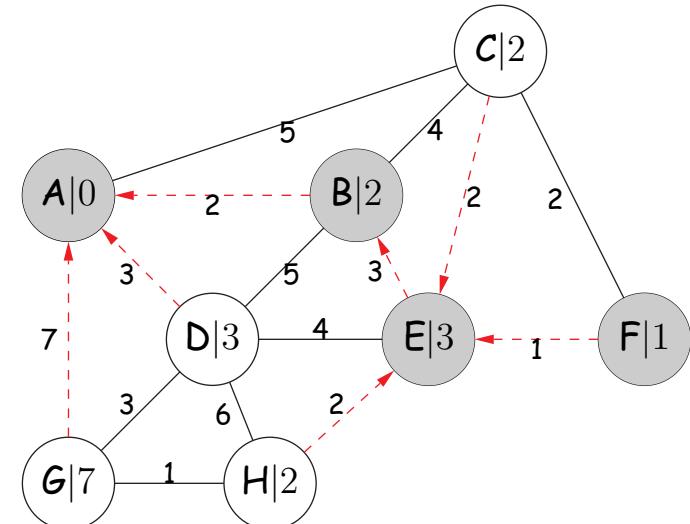
```
PriorityQueue fringe;  
For each node v { v.dist() = infinity; v.parent() = null; }  
Choose an arbitrary starting node, s;  
s.dist() = 0;  
fringe = priority queue ordered by smallest .dist();  
add all vertices to fringe;  
while (!fringe.isEmpty()) {  
    Vertex v = fringe.removeFirst();  
  
    For each edge(v,w) {  
        if (w ∈ fringe && weight(v,w) < w.dist())  
            { w.dist() = weight(v, w); w.parent() = v; }  
    }  
}
```



Minimum Spanning Trees by Prim's Algorithm

- Idea is to grow a tree starting from an arbitrary node.
- At each step, add the shortest edge connecting some node already in the tree to one that isn't yet.
- Why must this work?

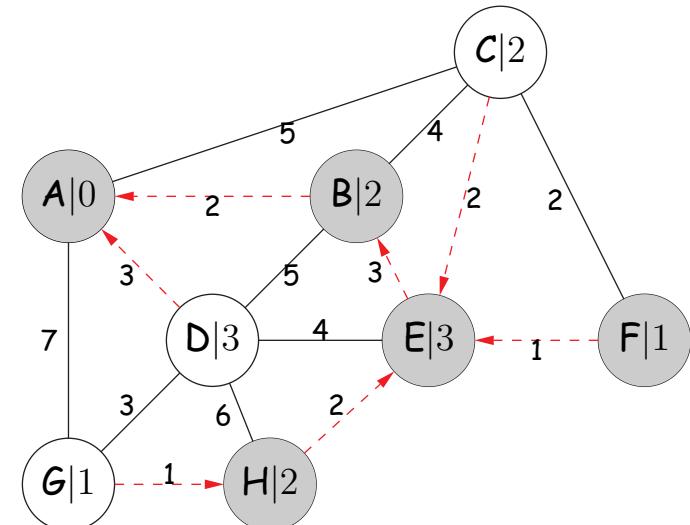
```
PriorityQueue fringe;  
For each node v { v.dist() = infinity; v.parent() = null; }  
Choose an arbitrary starting node, s;  
s.dist() = 0;  
fringe = priority queue ordered by smallest .dist();  
add all vertices to fringe;  
while (!fringe.isEmpty()) {  
    Vertex v = fringe.removeFirst();  
  
    For each edge(v,w) {  
        if (w ∈ fringe && weight(v,w) < w.dist())  
            { w.dist() = weight(v, w); w.parent() = v; }  
    }  
}
```



Minimum Spanning Trees by Prim's Algorithm

- Idea is to grow a tree starting from an arbitrary node.
- At each step, add the shortest edge connecting some node already in the tree to one that isn't yet.
- Why must this work?

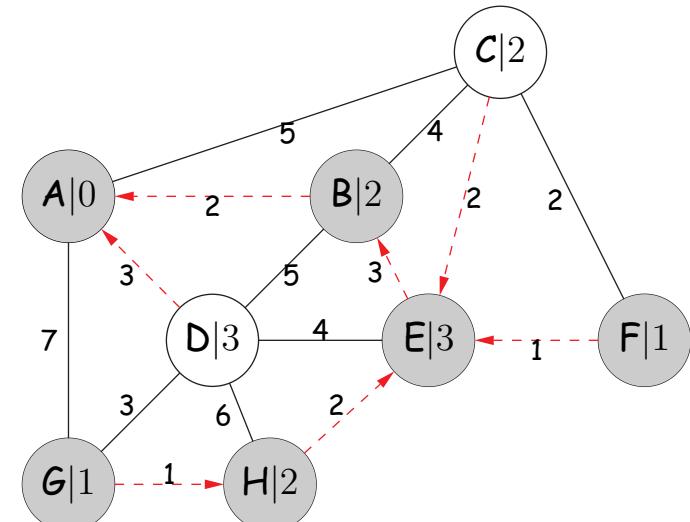
```
PriorityQueue fringe;  
For each node v { v.dist() = infinity; v.parent() = null; }  
Choose an arbitrary starting node, s;  
s.dist() = 0;  
fringe = priority queue ordered by smallest .dist();  
add all vertices to fringe;  
while (!fringe.isEmpty()) {  
    Vertex v = fringe.removeFirst();  
  
    For each edge(v,w) {  
        if (w ∈ fringe && weight(v,w) < w.dist())  
            { w.dist() = weight(v, w); w.parent() = v; }  
    }  
}
```



Minimum Spanning Trees by Prim's Algorithm

- Idea is to grow a tree starting from an arbitrary node.
- At each step, add the shortest edge connecting some node already in the tree to one that isn't yet.
- Why must this work?

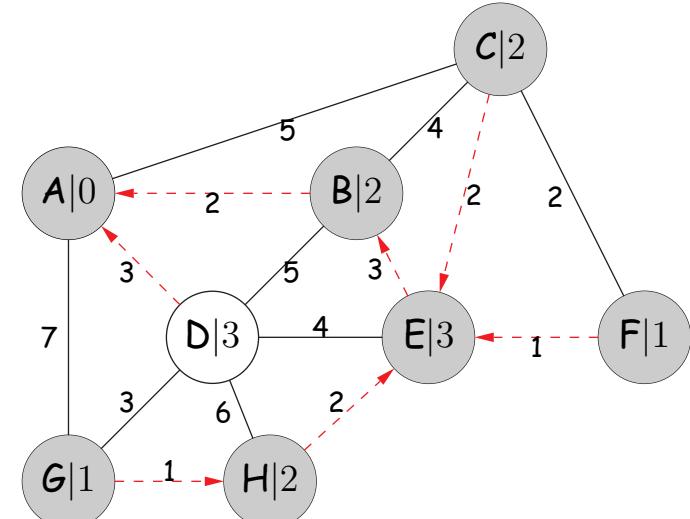
```
PriorityQueue fringe;  
For each node v { v.dist() = infinity; v.parent() = null; }  
Choose an arbitrary starting node, s;  
s.dist() = 0;  
fringe = priority queue ordered by smallest .dist();  
add all vertices to fringe;  
while (!fringe.isEmpty()) {  
    Vertex v = fringe.removeFirst();  
  
    For each edge(v,w) {  
        if (w ∈ fringe && weight(v,w) < w.dist())  
            { w.dist() = weight(v, w); w.parent() = v; }  
    }  
}
```



Minimum Spanning Trees by Prim's Algorithm

- Idea is to grow a tree starting from an arbitrary node.
- At each step, add the shortest edge connecting some node already in the tree to one that isn't yet.
- Why must this work?

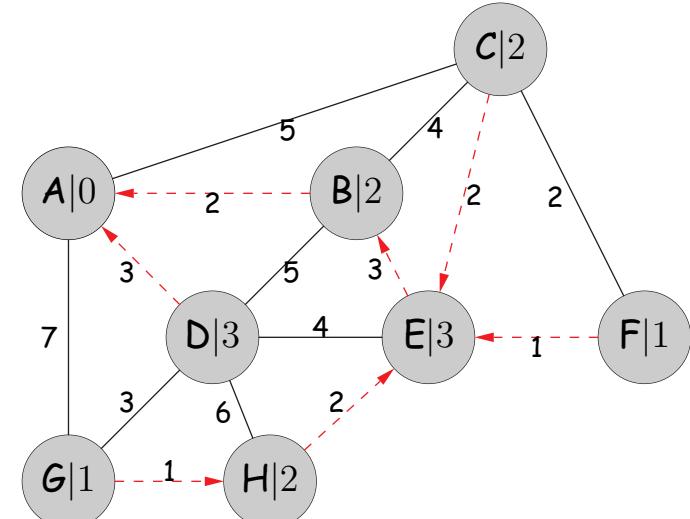
```
PriorityQueue fringe;  
For each node v { v.dist() = infinity; v.parent() = null; }  
Choose an arbitrary starting node, s;  
s.dist() = 0;  
fringe = priority queue ordered by smallest .dist();  
add all vertices to fringe;  
while (!fringe.isEmpty()) {  
    Vertex v = fringe.removeFirst();  
  
    For each edge(v,w) {  
        if (w ∈ fringe && weight(v,w) < w.dist())  
            { w.dist() = weight(v, w); w.parent() = v; }  
    }  
}
```



Minimum Spanning Trees by Prim's Algorithm

- Idea is to grow a tree starting from an arbitrary node.
- At each step, add the shortest edge connecting some node already in the tree to one that isn't yet.
- Why must this work?

```
PriorityQueue fringe;  
For each node v { v.dist() = ∞; v.parent() = null; }  
Choose an arbitrary starting node, s;  
s.dist() = 0;  
fringe = priority queue ordered by smallest .dist();  
add all vertices to fringe;  
while (!fringe.isEmpty()) {  
    Vertex v = fringe.removeFirst();  
  
    For each edge(v,w) {  
        if (w ∈ fringe && weight(v,w) < w.dist())  
            { w.dist() = weight(v, w); w.parent() = v; }  
    }  
}
```



Minimum Spanning Trees by Kruskal's Algorithm

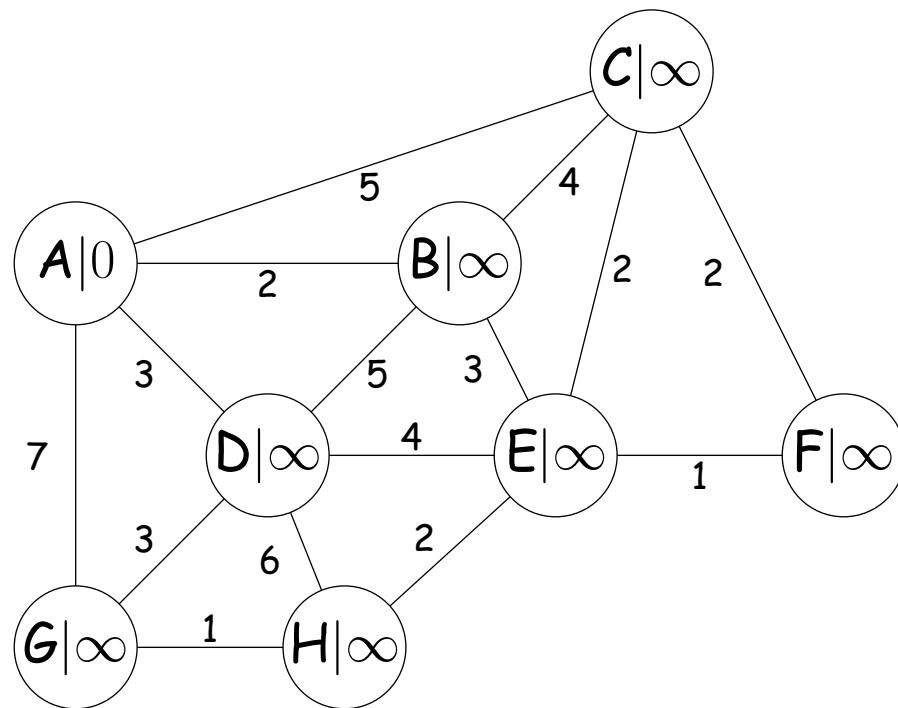
- Observation: the shortest edge in a graph can always be part of a minimum spanning tree.
- In fact, if we have a bunch of subtrees of a MST, then the shortest edge that connects two of them can be part of a MST, combining the two subtrees into a bigger one.
- So,...

Create one (trivial) subtree for each node in the graph;

MST = {};

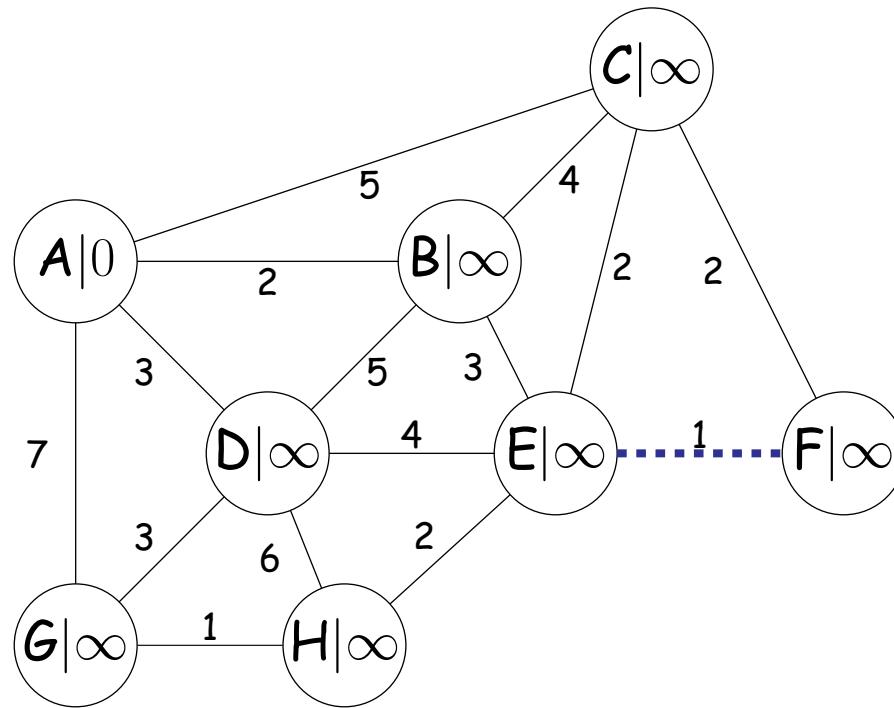
```
for each edge(v,w), in increasing order of weight {  
    if ( (v,w) connects two different subtrees ) {  
        Add (v,w) to MST;  
        Combine the two subtrees into one;  
    }  
}
```

Example of Kruskal's Algorithm



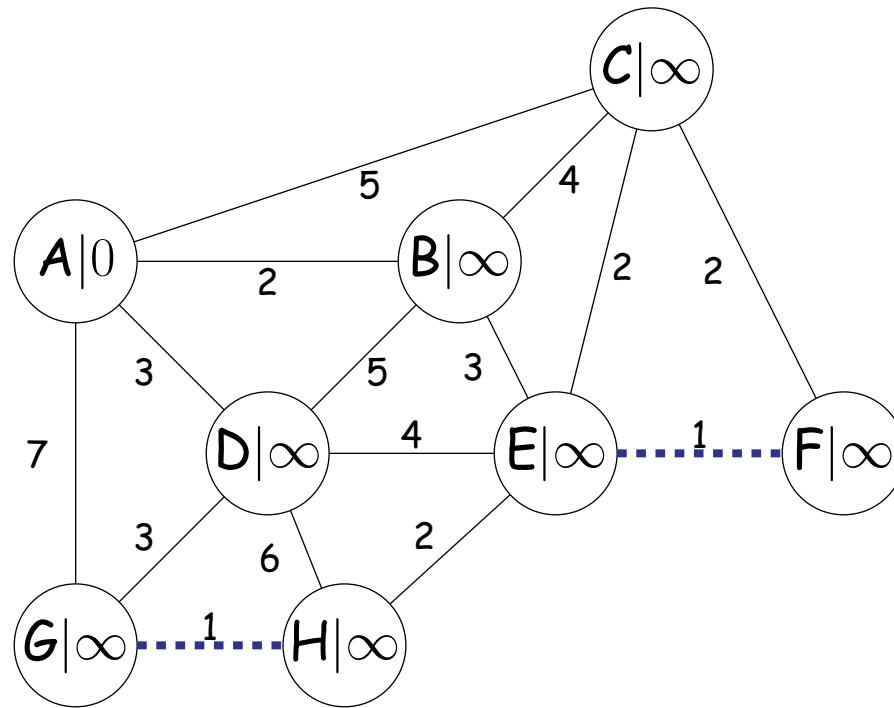
- We go through the edges in increasing order of weight (here, I break ties alphabetically).
- Edges that connect two unconnected groups get added to the tree (dashed lines).
- Edges that join already-joined groups are discarded ('X'ed out here).
- End result is a minimal spanning tree (a free tree).

Example of Kruskal's Algorithm



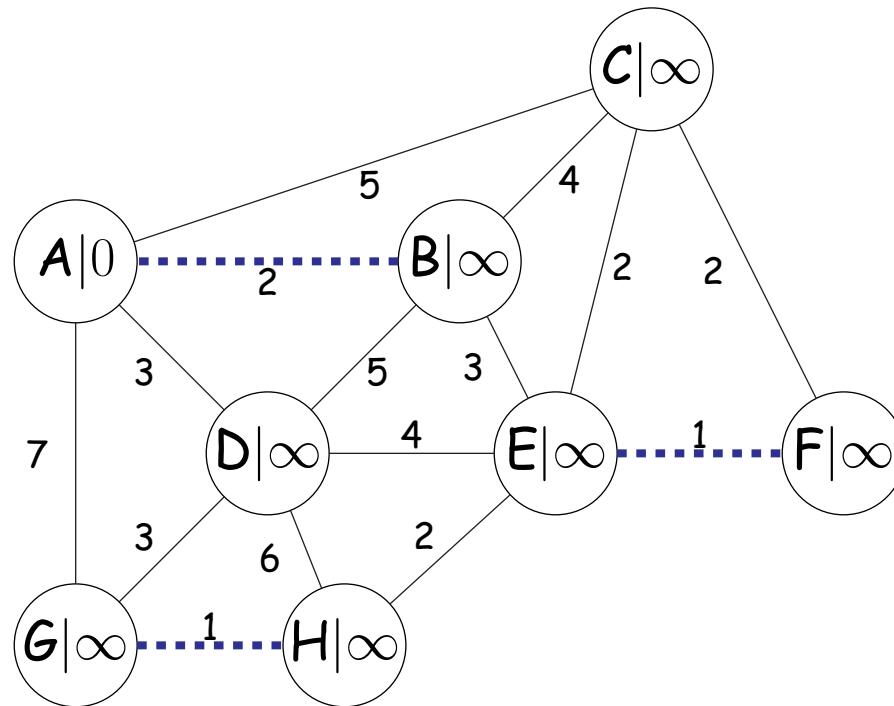
- We go through the edges in increasing order of weight (here, I break ties alphabetically).
- Edges that connect two unconnected groups get added to the tree (dashed lines).
- Edges that join already-joined groups are discarded ('X'ed out here).
- End result is a minimal spanning tree (a free tree).

Example of Kruskal's Algorithm



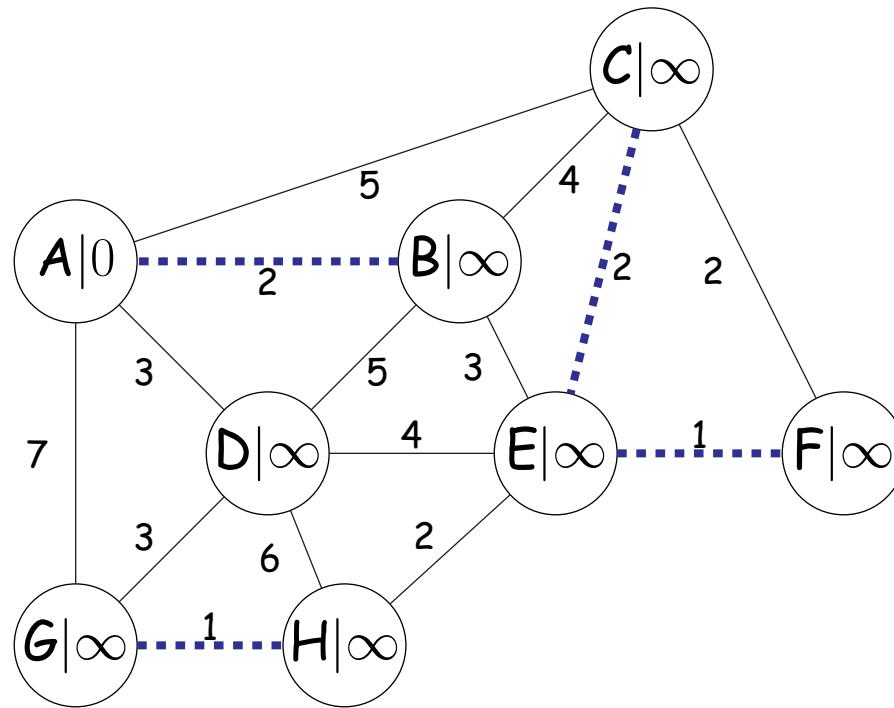
- We go through the edges in increasing order of weight (here, I break ties alphabetically).
- Edges that connect two unconnected groups get added to the tree (dashed lines).
- Edges that join already-joined groups are discarded ('X'ed out here).
- End result is a minimal spanning tree (a free tree).

Example of Kruskal's Algorithm



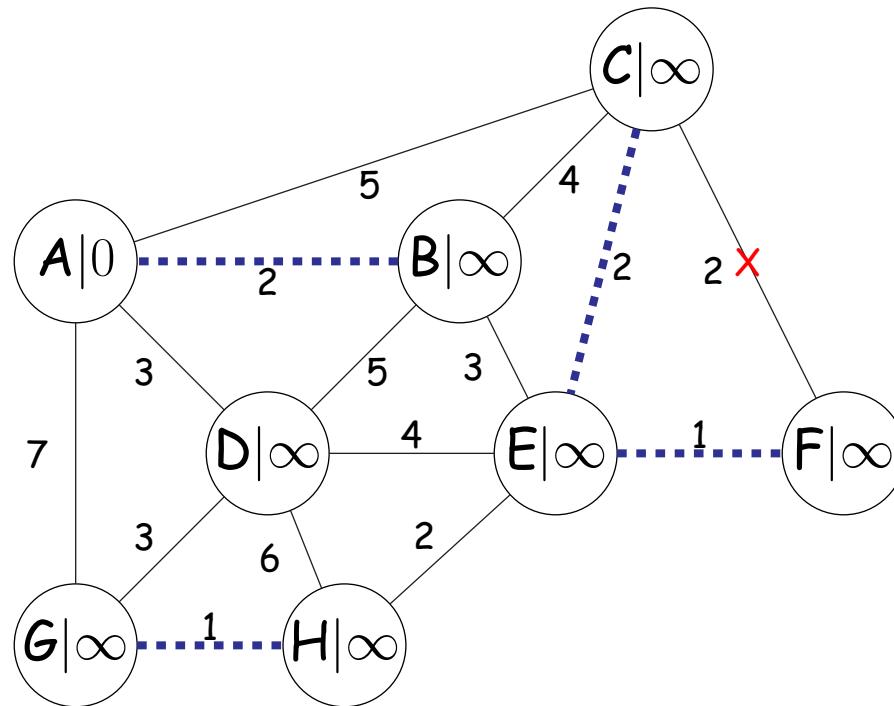
- We go through the edges in increasing order of weight (here, I break ties alphabetically).
- Edges that connect two unconnected groups get added to the tree (dashed lines).
- Edges that join already-joined groups are discarded ('X'ed out here).
- End result is a minimal spanning tree (a free tree).

Example of Kruskal's Algorithm



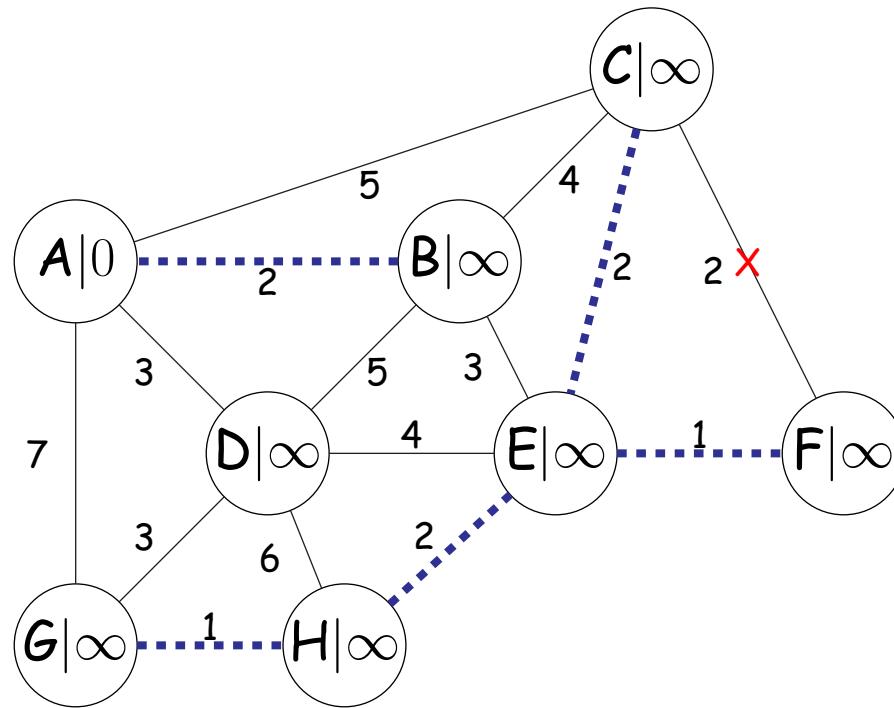
- We go through the edges in increasing order of weight (here, I break ties alphabetically).
- Edges that connect two unconnected groups get added to the tree (dashed lines).
- Edges that join already-joined groups are discarded ('X'ed out here).
- End result is a minimal spanning tree (a free tree).

Example of Kruskal's Algorithm



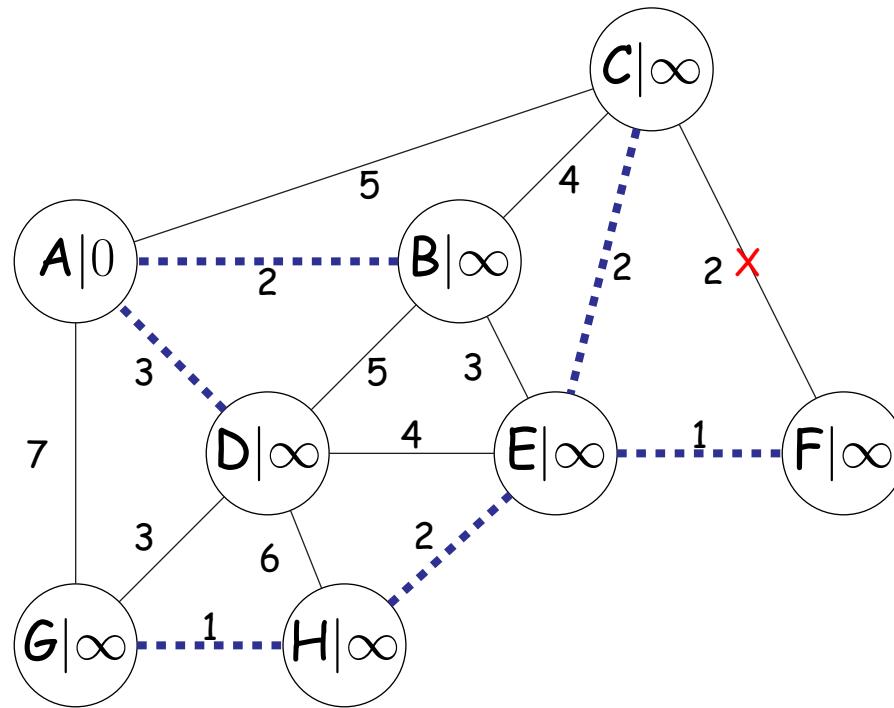
- We go through the edges in increasing order of weight (here, I break ties alphabetically).
- Edges that connect two unconnected groups get added to the tree (dashed lines).
- Edges that join already-joined groups are discarded ('X'ed out here).
- End result is a minimal spanning tree (a free tree).

Example of Kruskal's Algorithm



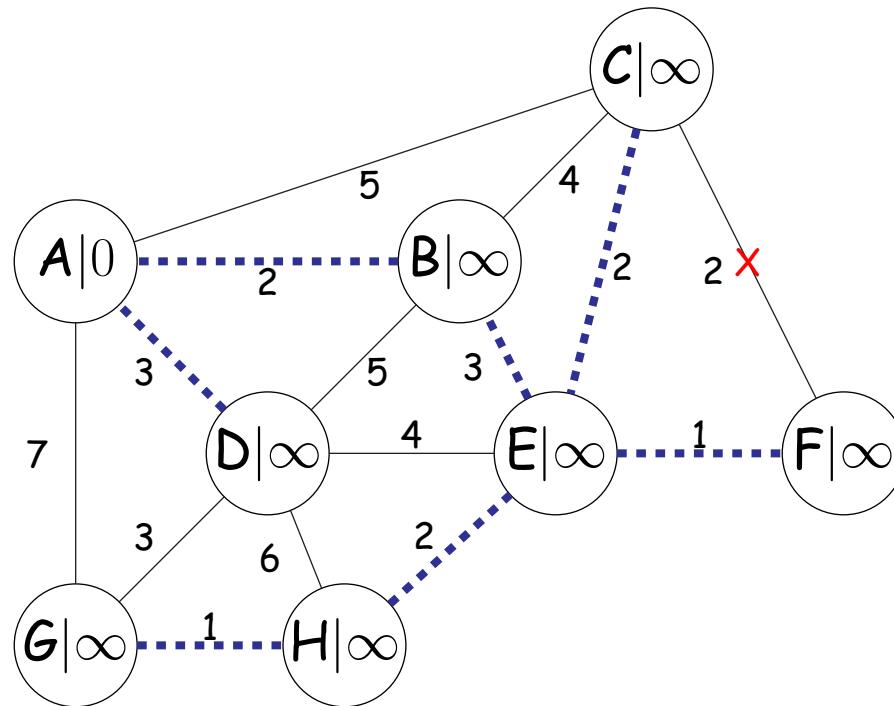
- We go through the edges in increasing order of weight (here, I break ties alphabetically).
- Edges that connect two unconnected groups get added to the tree (dashed lines).
- Edges that join already-joined groups are discarded ('X'ed out here).
- End result is a minimal spanning tree (a free tree).

Example of Kruskal's Algorithm



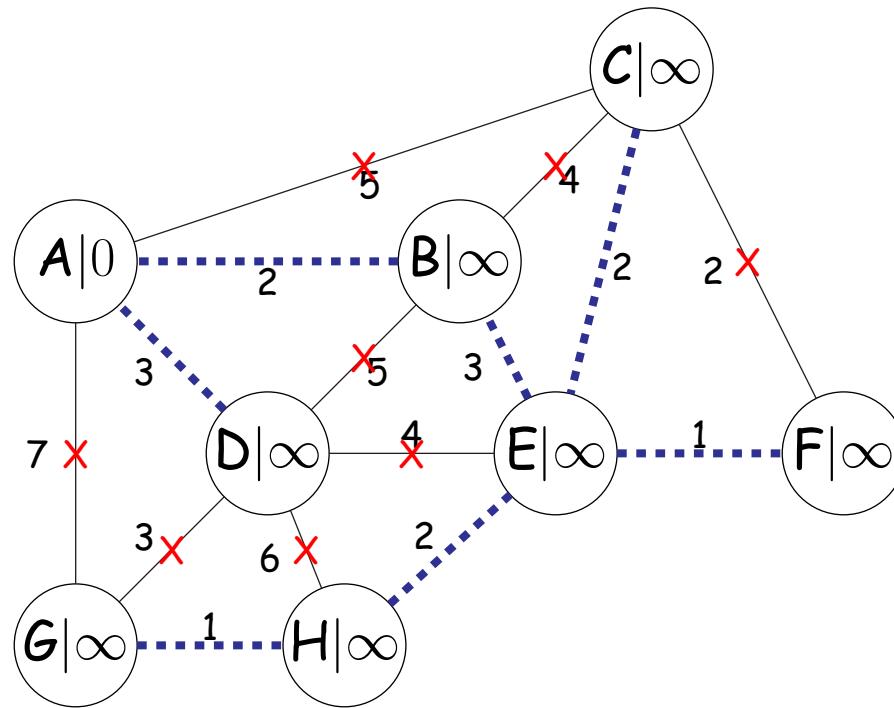
- We go through the edges in increasing order of weight (here, I break ties alphabetically).
- Edges that connect two unconnected groups get added to the tree (dashed lines).
- Edges that join already-joined groups are discarded ('X'ed out here).
- End result is a minimal spanning tree (a free tree).

Example of Kruskal's Algorithm



- We go through the edges in increasing order of weight (here, I break ties alphabetically).
- Edges that connect two unconnected groups get added to the tree (dashed lines).
- Edges that join already-joined groups are discarded ('X'ed out here).
- End result is a minimal spanning tree (a free tree).

Example of Kruskal's Algorithm



- We go through the edges in increasing order of weight (here, I break ties alphabetically).
- Edges that connect two unconnected groups get added to the tree (dashed lines).
- Edges that join already-joined groups are discarded ('X'ed out here).
- End result is a minimal spanning tree (a free tree).

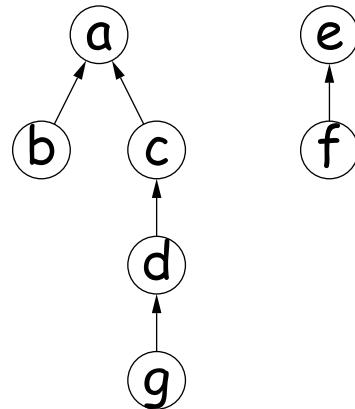
Union Find

- Kruskal's algorithm required that we have a set of sets of nodes with two operations:
 - *Find* which of the sets a given node belongs to.
 - Replace two sets with their *union*, reassigning all the nodes in the two original sets to this union.
- Obvious thing to do is to store a set number in each node, making finds fast.
- Union requires changing the set number in one of the two sets being merged; the smaller is better choice.
- This means an individual union can take $\Theta(N)$ time.
- Can union be faster?

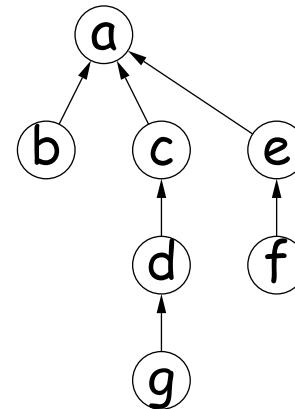
A Clever Trick

- Let's choose to represent a set of nodes by *one* arbitrary representative node in that set.
- Let every node contain a pointer to another node in the same set.
- Arrange for each pointer to represent the *parent* of a node in a tree that has the representative node as its root.
- To find what set a node is in, follow parent pointers.
- To union two such trees, make one root point to the other (choose the root of the larger tree as the union representative).

Two Sets

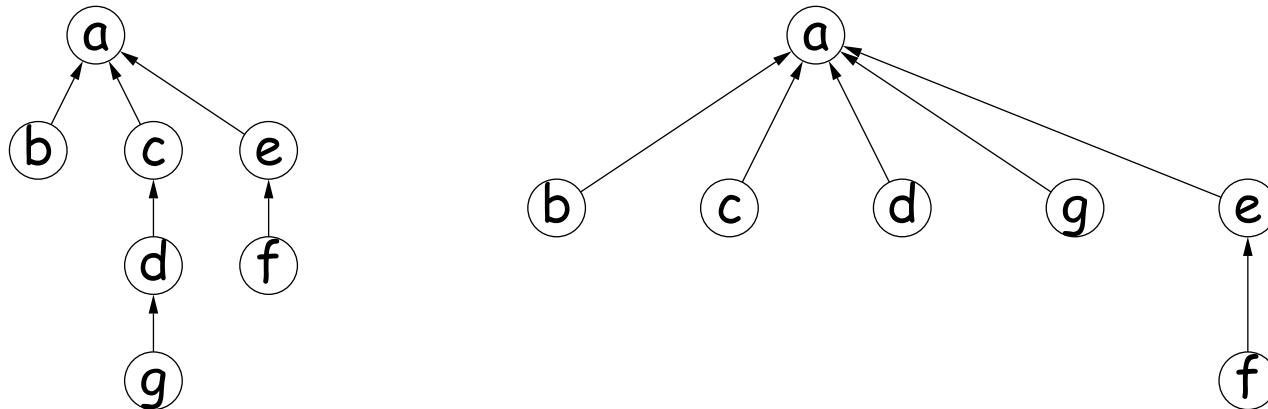


Their Union



Path Compression

- This makes unioning really fast, but the find operation is still ($\Omega(\lg N)$).
- So use the following trick: whenever we do a *find* operation, *compress* the path to the root, so that *subsequent* finds will be faster.
- That is, make each of the nodes in the path point directly to the root.
- Now union is very fast, and any sequence of unions and finds each have very, *very* nearly constant amortized time.
- Example: find 'g' in last tree (result of compression on right):



CS61B Lecture #35

Today:

- Pseudo-random Numbers (Chapter 11)
- What use are random sequences?
- What *are* "random sequences"?
- Pseudo-random sequences.
- How to get one.
- Relevant Java library classes and methods.
- Random permutations.

Why Random Sequences?

- Choose statistical samples
- Simulations
- Random algorithms
- Cryptography:
 - Choosing random keys and *nonces* (random one-time values used to make messages unique.)
 - Generating streams of random bits (e.g., stream ciphers encrypt messages by xor'ing reproducible streams of pseudo-random bits with the bits of the message.)
- And, of course, games

What Is a "Random Sequence"?

- How about: "a sequence where all numbers occur with equal frequency"?
 - Like 1, 2, 3, 4, ...?
- Well then, how about: "an unpredictable sequence where all numbers occur with equal frequency"?
 - Like 0, 0, 0, 1, 1, 2, 2, 2, 2, 2, 3, 4, 4, 0, 1, 1, 1, ...?
- Besides, what is wrong with 0, 0, 0, 0, ... anyway? Can't that occur by random selection?

Pseudo-Random Sequences

- Even if definable, a “truly” random sequence is difficult (i.e., slow) for a computer (or human) to produce. Must have some nondeterministic external source. Can use:
 - Periods between radioactive decays or cosmic rays.
 - Periods between keystrokes or incoming internet messages.
 - Mechanical coin flips or roulette wheels?
- For most purposes, we need only a sequence that satisfies certain statistical properties, even if deterministic (as is useful for reproducibility).
- Sometimes (e.g., cryptography) we need sequences that are *hard* or *impractical* to predict.
- *Pseudo-random sequence*: deterministic sequence that passes some given set of statistical tests that random sequences (probably) pass.
- For example, look at lengths of *runs*: increasing or decreasing contiguous subsequences.
- Unfortunately, statistical criteria to be used are quite involved. For details, see Knuth, volume 2.

Generating Pseudo-Random Sequences

- Not as easy as you might think: seemingly complex jumbling methods can give rise to bad sequences.
- *Linear congruential generators* are simple generators used by Java:

$X_0 = \text{arbitrary seed}$

$X_i = (aX_{i-1} + c) \bmod m, \quad i > 0$

- Usually, m is a power of 2.
- For best results, want $a \equiv 5 \pmod 8$, and a, c, m with no common factors.
- This gives generator with a *period of m* (length of sequence before repetition), and reasonable *potency* (a measure of certain dependencies among adjacent X_i .)
- Also want bits of a to "have no obvious pattern" and pass certain other tests (see Knuth).
- Java uses $a = 25214903917$, $c = 11$, $m = 2^{48}$, to compute 48-bit pseudo-random numbers. It's good enough for many purposes, but not *cryptographically secure*.

What Can Go Wrong (I)?

Any one who considers arithmetical methods of producing random digits is, of course, in a state of sin.

JOHN VON NEUMANN (1951)

- Short periods, many impossible values: E.g., a , c , m even.
- Obvious patterns. E.g., just using lower 3 bits of X_i in Java's 48-bit generator, to get integers in range 0 to 7. By properties of modular arithmetic,

$$\begin{aligned}X_i \bmod 8 &= (25214903917X_{i-1} + 11 \bmod 2^{48}) \bmod 8 \\&= (5(X_{i-1} \bmod 8) + 3) \bmod 8\end{aligned}$$

so we have a period of 8 on this generator; sequences like

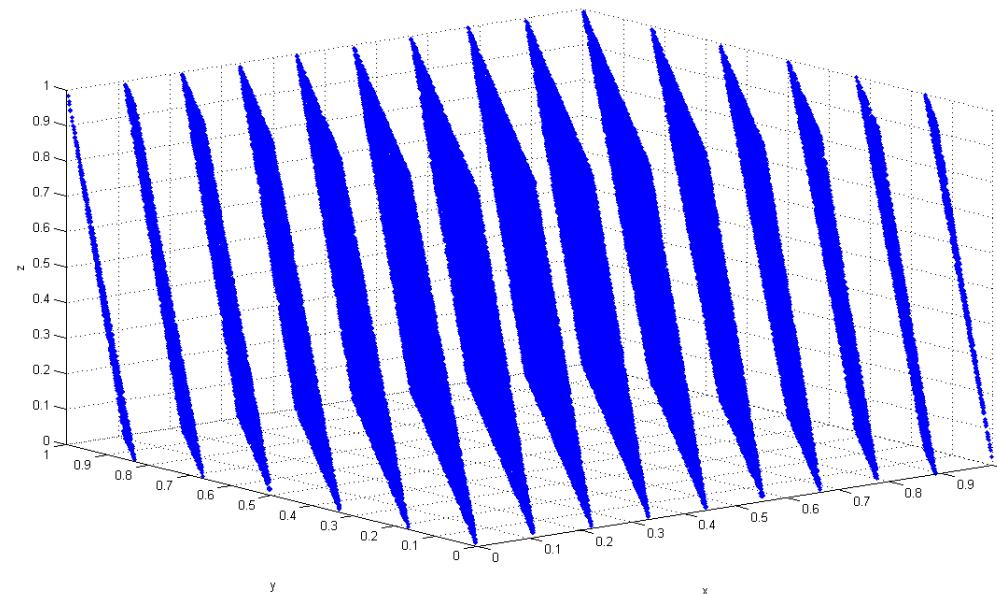
$$0, 1, 3, 7, 1, 2, 7, 1, 4, \dots$$

are impossible. This is why Java doesn't give you the raw 48 bits.

What Can Go Wrong (II)?

Bad potency leads to bad correlations.

- The infamous IBM generator RANDU: $c = 0$, $a = 65539$, $m = 2^{31}$.
- When RANDU is used to make 3D points: $(X_i/S, X_{i+1}/S, X_{i+2}/S)$, where S scales to a unit cube, ...
- ... points will be arranged in parallel planes with voids between. So "random points" won't ever get near many points in the cube:



[Credit: Luis Sanchez at English Wikipedia - Transferred from en.wikipedia to Commons by sevela.p., CC BY-SA 3.0, <https://commons.wikimedia.org/w/index.php?curid=3832343>]

Additive Generators

- Additive generator:

$$X_n = \begin{cases} \text{arbitrary value,} & n < 55 \\ (X_{n-24} + X_{n-55}) \bmod 2^e, & n \geq 55 \end{cases}$$

- Other choices than 24 and 55 possible.
- This one has period of $2^f(2^{55} - 1)$, for some $f < e$.
- Simple implementation with circular buffer:

```
i = (i+1) % 55;  
X[i] += X[(i+31) % 55]; // Why +31 (55-24) instead of -24?  
return X[i]; /* modulo 232 */
```

- where $X[0..54]$ is initialized to some “random” initial seed values.

Cryptographic Pseudo-Random Number Generators

- The simple form of linear congruential generators means that one can predict future values after seeing relatively few outputs.
- Not good if you want *unpredictable* output (think on-line games involving money or randomly generated keys for encrypting your web traffic.)
- A *cryptographic pseudo-random number generator (CPRNG)* has the properties that
 - Given k bits of a sequence, no polynomial-time algorithm can guess the next bit with better than 50% accuracy.
 - Given the current state of the generator, it is also infeasible to reconstruct the bits it generated in getting to that state.

Cryptographic Pseudo-Random Number Generator Example

- Start with a good *block cipher*—an encryption algorithm that encrypts blocks of N bits (not just one byte at a time as for Enigma). AES is an example.
- As a seed, provide a key, K , and an initialization value I .
- The j^{th} pseudo-random number is now $E(K, I + j)$, where $E(x, y)$ is the encryption of message y using key x .

Adjusting Range and Distribution

- Given a raw sequence of numbers, X_i , in range (e.g.) 0 to 2^{48} from the above methods, how do we get uniform random integers in range 0 to $n - 1$?
- If $n = 2^k$, it's easy: use the top k bits of next X_i (the bottom k bits are not as "random")
- For other n , be careful of slight biases at the ends. For example, if we compute X_i/M , where $M = \lfloor 2^{48}/n \rfloor$ (i.e., rounding down), then you can get n as a result (which you don't want: it's too big).
- If you try to fix that by setting $M = \lfloor 2^{48}/(n - 1) \rfloor$ instead, the probability of getting $n - 1$ will be wrong.
 - For example, suppose n is 300, so that $M = \lfloor 941387881975.4381 \dots \rfloor = 941387881975$.
 - There are M possible X values in the range $0..2^{48} - 1$ that give a result of 0, of 1, of 2, etc. But there are only 131 that give a result of 299 (the last number in the range.)

Adjusting Range (II)

- To fix this bias problems when n does not evenly divide 2^{48} , Java throws out values after the largest multiple of n that is less than 2^{48} :

```
/** Random integer in the range 0 .. n-1, n>0. */
int nextInt(int n) {
    long X = next random long ( $0 \leq X < 2^{48}$ );
    if (n is  $2^k$  for some k)
        return top k bits of X;

    int MAX = largest multiple of n that is < 248;
    while ( $X_i \geq MAX$ )
        X = next random long ( $0 \leq X < 2^{48}$ );
    return  $X_i / (MAX/n)$ ;
}
```

Arbitrary Bounds

- How to get arbitrary range of integers (L to U)?
- To get random float, x in range $0 \leq x < d$, compute

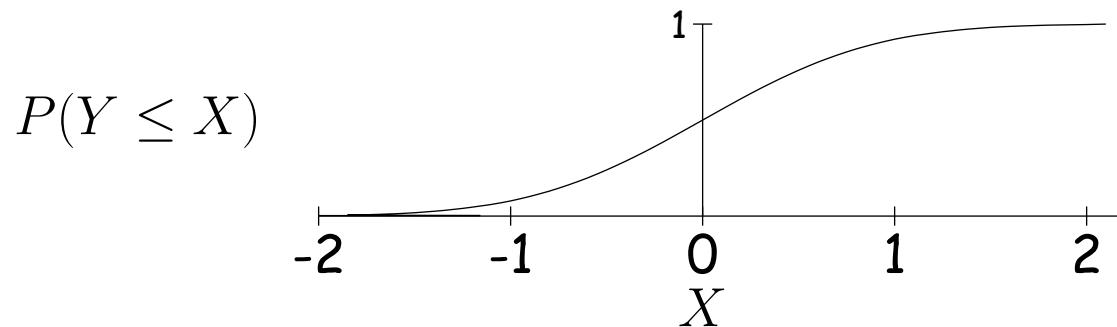
```
return d*nextInt(1<<24) / (1<<24);
```

- Random doubles are a bit more complicated: we need two integers to get enough bits.

```
long bigRand = ((long) nextInt(1<<26) << 27)
              + (long) nextInt(1<<27);
return d * bigRand / (1L << 53);
```

Generalizing: Other Distributions

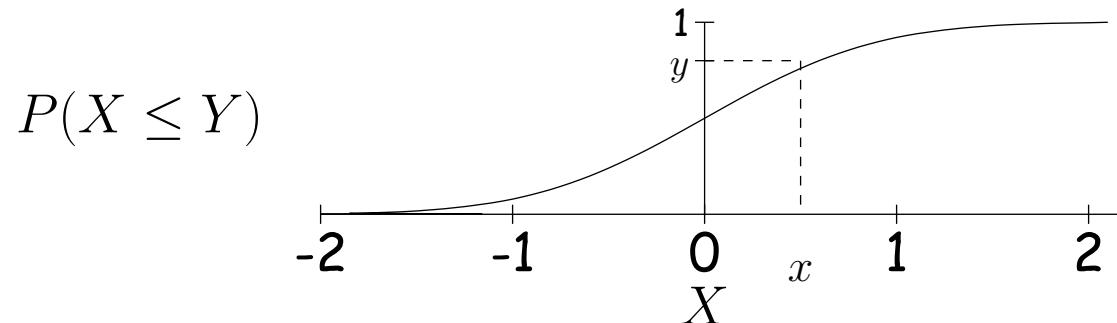
- Suppose we have some desired probability distribution function, and want to get random numbers that are distributed according to that distribution. How can we do this?
- Example: the normal distribution:



- Curve is the desired probability distribution. $P(Y \leq X)$ is the probability that random variable Y is $\leq X$.

Generalizing: Other Distributions (II)

Solution: Choose y uniformly between 0 and 1, and the corresponding x will be distributed according to P .



Java Classes

- `Math.random()`: random double in [0..1].
- Class `java.util.Random`: a random number generator with constructors:
 - `Random()` generator with “random” seed (based on time).
 - `Random(seed)` generator with given starting value (reproducible).
- Methods
 - `next(k)` k -bit random integer
 - `nextInt(n)` int in range [0.. n].
 - `nextLong()` random 64-bit integer.
 - `nextBoolean()`, `nextFloat()`, `nextDouble()` Next random values of other primitive types.
 - `nextGaussian()` normal distribution with mean 0 and standard deviation 1 (“bell curve”).
- `Collections.shuffle(L, R)` for list L and Random R permutes L randomly (using R).

Shuffling

- A *shuffle* is a random permutation of some sequence.
- Obvious dumb technique for shuffling an N -element list:
 - Generate N random numbers
 - Attach each to one of the list elements
 - Sort the list using the random numbers as keys.
- Can do quite a bit better:

```
void shuffle(List L, Random R) {  
    for (int i = L.size(); i > 0; i -= 1)  
        swap elements i-1 and R.nextInt(i) of L;  
}
```

- Example:

Swap items	0	1	2	3	4	5
Start	A♣	2♣	3♣	A♥	2♥	3♥
5 ↔ 1	A♣	3♥	3♣	A♥	2♥	2♣
4 ↔ 2	A♣	3♥	2♥	A♥	3♣	2♣

Swap items	0	1	2	3	4	5
3 ↔ 3	A♣	3♥	2♥	A♥	3♣	2♣
2 ↔ 0	2♥	3♥	A♣	A♥	3♣	2♣
1 ↔ 0	3♥	2♥	A♣	A♥	3♣	2♣

Random Selection

- The same technique would allow us to select N items from a list:

```
/** Permute L and return sublist of K>=0 randomly
 * chosen elements of L, using R as random source. */
List select(List L, int k, Random R) {
    for (int i = L.size(); i+k > L.size(); i -= 1)
        swap element i-1 of L with element
        R.nextInt(i) of L;
    return L.sublist(L.size()-k, L.size());
}
```

- Not terribly efficient for selecting random sequence of K distinct integers from $[0..N]$, with $K \ll N$.

Alternative Selection Algorithm (Floyd)

```
/** Random sequence of K distinct integers
 * from 0..N-1, 0<=K<=N. */
List<Integer> select(int N, int K, Random R)
{
    ArrayList<Integer> S = new ArrayList<>();

    for (int i = N-K; i < N; i += 1) {
        // All values in S are < i
        int s = R.nextInt(i+1); // 0 <= s <= i < N
        if (s == S.get(j) for some j)
            // Insert value i (which can't be there
            // yet) after the s (i.e., at a random
            // place other than the front)
            S.add(j+1, i);
        else
            // Insert random value s (which can't be
            // there yet) at front
            S.add(0, s);
    }
    return S;
}
```

Example

i	s	S
5	4	[4]
6	2	[2, 4]
7	5	[5, 2, 4]
8	5	[5, 8, 2, 4]
9	4	[5, 8, 2, 4, 9]

selectRandomIntegers(10, 5, R)

CS61B Lecture #36

Today:

- Dynamic Programming
- A Brief Side Trip: Enumeration types.

Dynamic Programming

- A puzzle (D. Garcia):
 - Start with a list with an even number of non-negative integers.
 - Each player in turn takes either the leftmost number or the rightmost.
 - Idea is to get the largest possible sum.
- Example: starting with (6, 12, 0, 8), you (as first player) should take the 8. Whatever the second player takes, you also get the 12, for a total of 20.
- Assuming your opponent plays perfectly (i.e., to get as much as possible), how can you maximize your sum?
- Can solve this with exhaustive game-tree search, but...

Obvious Program

- Recursion makes it easy, again:

```
int bestSum(int[] V) {
    int total, i, N = V.length;
    for (i = 0, total = 0; i < N; i += 1) total += V[i];
    return N == 0 ? 0 : bestSum(V, 0, N-1, total);
}

/** The largest sum obtainable by the first player in the game on
 * the list V[LEFT..RIGHT], given TOTAL = sum of V[LEFT..RIGHT]. */
int bestSum(int[] V, int left, int right, int total) {
    if (left == right) {
        return V[left];
    } else {
        int L = total - bestSum(V, left+1, right, total-V[left]);
        int R = total - bestSum(V, left, right-1, total-V[right]);
        return Math.max(L, R);
    }
}
```

- Time cost is $C(0) = 1$, $C(N) = 2C(N - 1)$; so $C(N) \in \Theta(2^N)$

Still Another Idea from CS61A

The problem is that we are recomputing intermediate results many times. Solution: *memoize* the intermediate results. Here, we pass in an $N \times N$ array ($N = V.length$) of memoized results, all initially -1.

```
int bestSum(int[] V) {
    int total, i, N = V.length;
    int[] memo = new int[N][N];
    for (i = 0, total = 0; i < N; i += 1) total += V[i];
    for (int[] row : memo) Arrays.fill(row, -1);
    return N == 0 ? 0 : bestSum(V, 0, N-1, total, memo);
}

int bestSum(int[] V, int left, int right, int total, int[][] memo) {
    if (left == right) memo[left][right] = V[left];
    else if (memo[left][right] == -1) {
        int L = total - bestSum(V, left+1, right, total-V[left], memo);
        int R = total - bestSum(V, left, right-1, total-V[right], memo);
        memo[left][right] = Math.max(L, R);
    }
    return memo[left][right];
}
```

Time Cost with Memoization

```
int bestSum(int[] V, int left, int right, int total, int[][] memo) {  
    if (left == right) memo[left][right] = V[left];  
    else if (memo[left][right] == -1) {  
        int L = total - bestSum(V, left+1, right, total-V[left], memo);  
        int R = total - bestSum(V, left, right-1, total-V[right], memo);  
        memo[left][right] = Math.max(L, R);  
    }  
    return memo[left][right];  
}
```

- Now we will perform the full computation of each entry in `memo` once.
- The number of recursive calls to `bestSum` must be $O(N^2)$, for $N =$ the length of V , an enormous improvement from $\Theta(2^N)$!

Iterative Version

- In the usual presentation of this idea—known as *dynamic programming*—we figure out ahead of time the order in which the memoized version will fill in `memo`, and write an explicit loop, saving the time needed to check whether the result has already been computed.
- Suppose that we use as input $\{ 2 \ 1 \ 8 \ 3 \}$. Our goal is to find the circled—(0,3)—entry in the memo table, which gets filled in like this (grayed boxes indicates table entries that are used to compute a new entry):

		right			
		0	1	2	3
		0			
		1			
		2			
		3			

Iterative Version

- In the usual presentation of this idea—known as *dynamic programming*—we figure out ahead of time the order in which the memoized version will fill in `memo`, and write an explicit loop, saving the time needed to check whether the result has already been computed.
- Suppose that we use as input $\{ 2 \ 1 \ 8 \ 3 \}$. Our goal is to find the circled—(0,3)—entry in the memo table, which gets filled in like this (grayed boxes indicates table entries that are used to compute a new entry):

		right			
		0	1	2	3
		0			
		1			
		2			
		3			3

Iterative Version

- In the usual presentation of this idea—known as *dynamic programming*—we figure out ahead of time the order in which the memoized version will fill in `memo`, and write an explicit loop, saving the time needed to check whether the result has already been computed.
- Suppose that we use as input $\{ 2 \ 1 \ 8 \ 3 \}$. Our goal is to find the circled—(0,3)—entry in the memo table, which gets filled in like this (grayed boxes indicates table entries that are used to compute a new entry):

		right			
		0	1	2	3
		0			
		1			
		2		8	
		3			3

Iterative Version

- In the usual presentation of this idea—known as *dynamic programming*—we figure out ahead of time the order in which the memoized version will fill in `memo`, and write an explicit loop, saving the time needed to check whether the result has already been computed.
- Suppose that we use as input $\{ 2 \ 1 \ 8 \ 3 \}$. Our goal is to find the circled—(0,3)—entry in the memo table, which gets filled in like this (grayed boxes indicates table entries that are used to compute a new entry):

		right			
		0	1	2	3
		0			
		1			
		2		8	8
		3			3

Iterative Version

- In the usual presentation of this idea—known as *dynamic programming*—we figure out ahead of time the order in which the memoized version will fill in `memo`, and write an explicit loop, saving the time needed to check whether the result has already been computed.
- Suppose that we use as input $\{ 2 \ 1 \ 8 \ 3 \}$. Our goal is to find the circled—(0,3)—entry in the memo table, which gets filled in like this (grayed boxes indicates table entries that are used to compute a new entry):

		right			
		0	1	2	3
		0			
		1	1		
		2		8	8
		3			3

Iterative Version

- In the usual presentation of this idea—known as *dynamic programming*—we figure out ahead of time the order in which the memoized version will fill in `memo`, and write an explicit loop, saving the time needed to check whether the result has already been computed.
- Suppose that we use as input $\{ 2 \ 1 \ 8 \ 3 \}$. Our goal is to find the circled—(0,3)—entry in the memo table, which gets filled in like this (grayed boxes indicates table entries that are used to compute a new entry):

		right			
		0	1	2	3
		0			
left		1	8		
2			8	8	
3					3

Iterative Version

- In the usual presentation of this idea—known as *dynamic programming*—we figure out ahead of time the order in which the memoized version will fill in `memo`, and write an explicit loop, saving the time needed to check whether the result has already been computed.
- Suppose that we use as input $\{ 2 \ 1 \ 8 \ 3 \}$. Our goal is to find the circled—(0,3)—entry in the memo table, which gets filled in like this (grayed boxes indicates table entries that are used to compute a new entry):

		right			
		0	1	2	3
		0			
left		1	8	4	
2			8	8	
3					3

Iterative Version

- In the usual presentation of this idea—known as *dynamic programming*—we figure out ahead of time the order in which the memoized version will fill in `memo`, and write an explicit loop, saving the time needed to check whether the result has already been computed.
- Suppose that we use as input $\{ 2 \ 1 \ 8 \ 3 \}$. Our goal is to find the circled—(0,3)—entry in the memo table, which gets filled in like this (grayed boxes indicates table entries that are used to compute a new entry):

		right			
		0	1	2	3
		2			
			1	8	4
				8	8
					3

Iterative Version

- In the usual presentation of this idea—known as *dynamic programming*—we figure out ahead of time the order in which the memoized version will fill in `memo`, and write an explicit loop, saving the time needed to check whether the result has already been computed.
- Suppose that we use as input $\{ 2 \ 1 \ 8 \ 3 \}$. Our goal is to find the circled—(0,3)—entry in the memo table, which gets filled in like this (grayed boxes indicates table entries that are used to compute a new entry):

		right				
		0	1	2	3	
		0	2	2		
		1		1	8	4
		2			8	8
		3				3

Iterative Version

- In the usual presentation of this idea—known as *dynamic programming*—we figure out ahead of time the order in which the memoized version will fill in `memo`, and write an explicit loop, saving the time needed to check whether the result has already been computed.
- Suppose that we use as input $\{ 2 \ 1 \ 8 \ 3 \}$. Our goal is to find the circled—(0,3)—entry in the memo table, which gets filled in like this (grayed boxes indicates table entries that are used to compute a new entry):

		right				
		0	1	2	3	
		0	2	2	9	3
		1		1	8	4
		2			8	8
		3				3

Iterative Version

- In the usual presentation of this idea—known as *dynamic programming*—we figure out ahead of time the order in which the memoized version will fill in `memo`, and write an explicit loop, saving the time needed to check whether the result has already been computed.
- Suppose that we use as input $\{ 2 \ 1 \ 8 \ 3 \}$. Our goal is to find the circled—(0,3)—entry in the memo table, which gets filled in like this (grayed boxes indicates table entries that are used to compute a new entry):

		right				
		0	1	2	3	
		0	2	2	9	10
		1		1	8	4
		2			8	8
		3				3

The Iterative Program

- To fill in the table in this order, we can use the following:

```
int bestSum(int[] V) {  
    int[][] memo = new int[V.length][V.length];  
    int total;  
    for (int left = V.length - 1; left >= 0; left -= 1) {  
        memo[left][left] = total = V[left];  
        for (int right = left + 1; right < V.length; right += 1) {  
            total += V[right];  
            int L = total - memo[left+1][right];  
            int R = total - memo[left][right-1];  
            memo[left][right] = Math.max(L, R);  
        }  
    }  
    return memo[0][V.length-1];  
}
```

- Generally, though, I say why bother unless the recursion would be too deep?

Compressing the Memo Table

- In fact, we can do even better, and use a one-dimensional memo table with only N entries.
- If you look at the previous animation of filling in the memo table, you'll see that we only ever reference entries immediately below and to the left of the item being computed.
- So we can actually just keep updating the same row.

```
int bestSum(int[] V) {  
    int[] memo = new int[V.length];  
    int total;  
    for (int left = V.length - 1; left >= 0; left -= 1) {  
        memo[left] = total = V[left];  
        for (int right = left + 1; right < V.length; right += 1) {  
            total += V[right];  
            int L = total - memo[right];  
            int R = total - memo[right-1];  
            memo[right] = Math.max(L, R);  
        }  
    }  
    return memo[V.length-1];  
}
```

Longest Common Subsequence

- **Problem:** Find length of the longest string that is a subsequence of each of two other strings.
- **Example:** Longest common subsequence of "sally\u2022sells\u2022sea\u2022shells\u2022by\u2022the\u2022seashore" and "sarah\u2022sold\u2022salt\u2022sellers\u2022at\u2022the\u2022salt\u2022mines" is "sa\u2022sl\u2022sa\u2022sells\u2022\u2022the\u2022sae" (length 23)
- Similarity testing, for example.
- Obvious recursive algorithm:

```
/** Length of longest common subsequence of S0[0..k0-1]
 * and S1[0..k1-1] (pseudo Java) */
static int lls(String S0, int k0, String S1, int k1) {
    if (k0 == 0 || k1 == 0) return 0;
    if (S0[k0-1] == S1[k1-1]) return 1 + lls(S0, k0-1, S1, k1-1);
    else return Math.max(lls(S0, k0-1, S1, k1), lls(S0, k0, S1, k1-1));
}
```

- Exponential, but obviously memoizable.

Memoized Longest Common Subsequence

```
/** Length of longest common subsequence of S0[0..k0-1] and S1[0..k1-1]. */
static int lls(String S0, int k0, String S1, int k1) {
    int[][] memo = new int[k0+1][k1+1];
    for (int[] row : memo) Arrays.fill(row, -1);
    return lls(S0, k0, S1, k1, memo);
}
private static int lls(String S0, int k0, String S1, int k1, int[][] memo) {
    if (k0 == 0 || k1 == 0) return 0;
    if (memo[k0][k1] == -1) {
        if (S0[k0-1] == S1[k1-1])
            memo[k0][k1] = 1 + lls(S0, k0-1, S1, k1-1, memo);
        else
            memo[k0][k1] = Math.max(lls(S0, k0-1, S1, k1, memo),
                                   lls(S0, k0, S1, k1-1, memo));
    }
    return memo[k0][k1];
}
```

Q: How fast will the memoized version be?

Memoized Longest Common Subsequence

```
/** Length of longest common subsequence of S0[0..k0-1] and S1[0..k1-1]. */
static int lls(String S0, int k0, String S1, int k1) {
    int[][] memo = new int[k0+1][k1+1];
    for (int[] row : memo) Arrays.fill(row, -1);
    return lls(S0, k0, S1, k1, memo);
}
private static int lls(String S0, int k0, String S1, int k1, int[][] memo) {
    if (k0 == 0 || k1 == 0) return 0;
    if (memo[k0][k1] == -1) {
        if (S0[k0-1] == S1[k1-1])
            memo[k0][k1] = 1 + lls(S0, k0-1, S1, k1-1, memo);
        else
            memo[k0][k1] = Math.max(lls(S0, k0-1, S1, k1, memo),
                                   lls(S0, k0, S1, k1-1, memo));
    }
    return memo[k0][k1];
}
```

Q: How fast will the memoized version be? $\Theta(k_0 \cdot k_1)$

Side Trip into Java: Enumeration Types

- Problem: Need a type with a few, named, discrete values.
- In the purest form, the only necessary operations are == and !=; the only property of a value of the type is that it differs from all others.
- In older versions of Java, used named integer constants:

```
interface Pieces {  
    int BLACK_PIECE = 0,      // Fields in interfaces are static final.  
    BLACK_KING = 1,  
    WHITE_PIECE = 2,  
    WHITE_KING = 3,  
    EMPTY = 4;  
}
```

- C and C++ provide *enumeration types* as a shorthand, with syntax like this:

```
enum Piece { BLACK_PIECE, BLACK_KING, WHITE_PIECE, WHITE_KING, EMPTY };
```

- But since all these values are basically ints, accidents can happen.

Enum Types in Java

- More modern versions of Java allow syntax like that of C or C++, but with more guarantees:

```
public enum Piece
    { BLACK_PIECE, BLACK_KING, WHITE_PIECE, WHITE_KING, EMPTY; }
```

- Defines Piece as a new reference type, a special kind of class type.
- The names BLACK_PIECE, etc., are static, final *enumeration constants* (or *enumerals*) of type PIECE.
- They are automatically initialized, and are the only values of the enumeration type that exist (cannot say `new Piece()`).
- Can safely use `==`, and also switch statements:

```
boolean isKing(Piece p) {
    switch (p) {
        case BLACK_KING: case WHITE_KING: return true;
        default: return false;
    }
}
```

Making Enumerals Available Elsewhere

- Enumerals like BLACK_PIECE are static members of a class, not classes.
- Therefore, unlike C or C++, their declarations are not automatically visible outside the enumeration class definition.
- So, in other classes, must write Piece.BLACK_PIECE, which can get annoying.
- However, with version 1.5, Java has *static imports*: to import all static definitions of class checkers.Piece (including enumerals), you write

```
import static checkers.Piece.*;
```

among the import clauses.

- Alas, *cannot* use this for enum classes in the anonymous package.

Operations on Enum Types

- Order of declaration of enumeration constants significant: `.ordinal()` gives the position (numbering from 0) of an enumeration value. Thus, `Piece.BLACK_KING.ordinal() is 1.`
- The array `Piece.values()` gives all the possible values of the type. Thus, you can write:

```
for (Piece p : Piece.values())
    System.out.printf("Piece value # %d is %s%n", p.ordinal(), p);
```

- The static function `Piece.valueOf` converts a String into a value of type `Piece`. So `Piece.valueOf("EMPTY") == EMPTY`.

Fancy Enum Types

- Enums are classes. You can define all the extra fields, methods, and constructors you want.
- Constructors are used only in creating enumeration constants. The constructor arguments follow the constant name:

```
enum Piece {  
    BLACK_PIECE(BLACK, false, "b"), BLACK_KING(BLACK, true, "B"),  
    WHITE_PIECE(WHITE, false, "w"), WHITE_KING(WHITE, true, "W"),  
    EMPTY(null, false, " ");  
  
    private final Side color;  
    private final boolean isKing;  
    private final String textName;  
  
    Piece(Side color, boolean isKing, String textName) {  
        this.color = color; this.isKing = isKing; this.textName = textName;  
    }  
    Side color() { return color; }  
    boolean isKing() { return isKing; }  
    String textName() { return textName; }  
}
```

Lecture #37

Today: Side excursion into nitty-gritty stuff: Threads.

Threads

- So far, all our programs consist of single sequences of instructions.
- Each such sequence is called a *thread* (for “thread of control”) in Java.
- Java supports programs containing *multiple* threads, which (conceptually) run concurrently.
- To allow program access to threads, Java provides the type `Thread` in `java.lang`. Each `Thread` contains information about, and controls, one thread.
- Simultaneous access to data from two threads can cause chaos, so there are also constructs for controlled communication, allowing threads to *lock* objects, to *wait* to be notified of events, and to *interrupt* other threads.

But Why?

- Actually, on a uniprocessor, only one thread at a time actually runs, while others wait, but this is largely invisible. So why bother with threads?
- Typical Java programs always have > 1 thread: besides the main program, others clean up garbage objects, receive signals, update the display, other stuff.
- When programs deal with asynchronous events, it is sometimes convenient to organize into subprograms, one for each independent, related sequence of events.
- Threads allow us to insulate one such subprogram from another. They thus provide a form of modularization.
- GUIs often organized like this: application is doing some computation or I/O, another thread waits for mouse clicks (like 'Stop'), another pays attention to updating the screen as needed.
- Large servers like search engines may be organized this way, with one thread per request.
- And, of course, sometimes we *do* have a real multiprocessor.

Java Mechanics

- To specify the actions "walking" and "chewing gum":

```
class Chewer1 implements Runnable {  
    public void run()  
    { while (true) ChewGum(); }  
}  
  
class Walker1 implements Runnable {  
    public void run()  
    { while (true) Walk(); }  
}
```

```
// Walk and chew gum  
Thread chomp  
= new Thread(new Chewer1());  
Thread clomp  
= new Thread(new Walker1());  
chomp.start(); clomp.start();
```

- Concise Alternative (uses fact that Thread implements Runnable):

```
class Chewer2 extends Thread {  
    public void run()  
    { while (true) ChewGum(); }  
}  
  
class Walker2 extends Thread {  
    public void run()  
    { while (true) Walk(); }  
}
```

```
Thread chomp = new Chewer2(),  
      clomp = new Walker2();  
chomp.start();  
clomp.start();
```

Avoiding Interference

- When one thread has data for another, one must wait for the other to be ready.
- Likewise, if two threads use the same object, generally only one should modify it at a time; other must wait.
- E.g., what would happen if two threads simultaneously inserted an item into a linked list at the same point in the list?
- A: Both could conceivably execute

```
p.next = new ListCell(x, p.next);
```

with the *same* values of p and p.next; one insertion is lost.

- Can arrange for only one thread at a time to execute a method on a particular object with either of the following equivalent definitions:

```
void f(...) {  
    synchronized (this) {  
        body of f  
    }  
}
```

```
synchronized void f(...) {  
    body of f  
}
```

Communicating the Hard Way

- Communicating data is tricky: the faster party must wait for the slower.
- Obvious approaches for sending data from thread to thread don't work:

```
class DataExchanger {  
    Object value = null;  
    Object receive() {  
        Object r; r = null;  
        while (r == null)  
            { r = value; }  
        value = null;  
        return r;  
    }  
    void deposit(Object data) {  
        while (value != null) { }  
        value = data;  
    }  
}
```

```
DataExchanger exchanger  
= new DataExchanger();  
-----  
// thread1 sends to thread2 with  
exchanger.deposit("Hello!");  
-----  
// thread2 receives from thread1 with  
msg = (String) exchanger.receive();
```

- **BAD:** One thread can monopolize the machine while waiting; two threads executing deposit or receive simultaneously cause chaos.

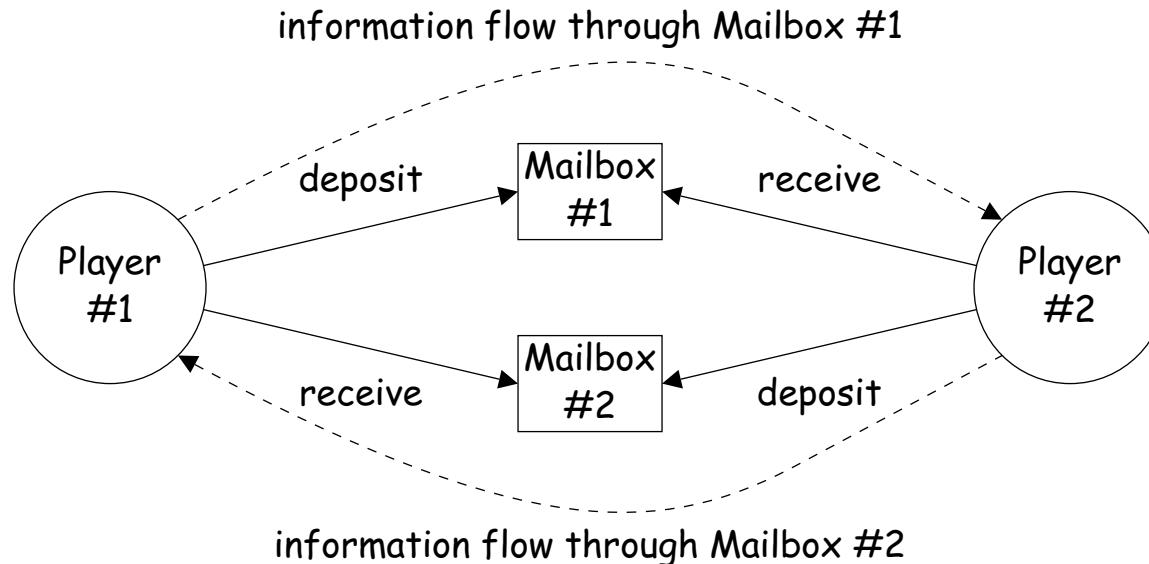
Primitive Java Facilities

- `Object.wait` method makes the current thread wait (not using processor) until notified by `notifyAll`, unlocking the Object while it waits.
- For example, `ucb.util.mailbox` has something like this (simplified):

```
interface Mailbox {  
    void deposit(Object msg) throws InterruptedException;  
    Object receive() throws InterruptedException;  
}  
  
class QueuedMailbox implements Mailbox {  
    private List<Object> queue = new LinkedList<Object>();  
  
    public synchronized void deposit(Object msg) {  
        queue.add(msg);  
        this.notifyAll(); // Wake any waiting receivers  
    }  
  
    public synchronized Object receive() throws InterruptedException {  
        while (queue.isEmpty()) wait();  
        return queue.remove(0);  
    }  
}
```

Message-Passing Style

- The use of Java primitives very error-prone. CS162 goes into alternatives.
- Mailboxes are higher-level, and allow the following program structure:



- Where each Player is a thread that looks like this:

```
while (! gameOver()) {  
    if (myMove())  
        outBox.deposit(computeMyMove(lastMove));  
    else  
        lastMove = inBox.receive();  
}
```

More Concurrency

- Previous example can be done other ways, but mechanism is very flexible.
- E.g., suppose you want to think during opponent's move:

```
while (!gameOver()) {  
    if (myMove())  
        outBox.deposit(computeMyMove(lastMove));  
    else {  
        do {  
            thinkAheadALittle();  
            lastMove = inBox.receiveIfPossible();  
        } while (lastMove == null);  
    }  
}
```

- `receiveIfPossible` (written `receive(0)` in our actual package) doesn't wait; returns null if no message yet, perhaps like this:

```
public synchronized Object receiveIfPossible()  
    throws InterruptedException {  
    if (queue.isEmpty())  
        return null;  
    return queue.remove(0);  
}
```

Coroutines

- A *coroutine* is a kind of synchronous thread that explicitly hands off control to other coroutines so that only one executes at a time, like Python generators. Can get similar effect with threads and mailboxes.
- Example: recursive inorder tree iterator:

```
class TreeIterator extends Thread {  
    Tree root; Mailbox r;  
    TreeIterator(Tree T, Mailbox r) {  
        this.root = T; this.dest = r;    void treeProcessor(Tree T) {  
    }                                Mailbox m = new QueuedMailbox();  
    public void run() {  
        traverse(root);  
        r.deposit(End marker);  
    }  
    void traverse(Tree t) {  
        if (t == null) return;  
        traverse(t.left);  
        r.deposit(t.label);  
        traverse(t.right);  
    }  
}  
Last modified: Mon Apr 25 11:29:29 2022
```

```
        new TreeIterator(T, m).start();  
        while (true) {  
            Object x = m.receive();  
            if (x is end marker)  
                break;  
            do something with x;  
        }  
    }
```

Use In GUIs

- Java runtime library uses a special thread that does nothing but wait for *events* like mouse clicks, pressed keys, mouse movement, etc.
- You can designate an object of your choice as a *listener*; which means that Java's event thread calls a method of that object whenever an event occurs.
- As a result, your program can do work while the UI continues to respond to buttons, menus, etc.
- Another special thread does all the drawing. You don't have to be aware when this takes place; just ask that the thread wake up whenever you change something.

Highlights of a GUI Component

```
/** A widget that draws multi-colored lines indicated by mouse. */
class Lines extends JComponent implements MouseListener {
    private List<Point> lines = new ArrayList<Point>();

    Lines() { // Main thread calls this to create one
        setPreferredSize(new Dimension(400, 400));
        addMouseListener(this);
    }
    public synchronized void paintComponent(Graphics g) { // Paint thread
        g.setColor(Color.white);    g.fillRect(0, 0, 400, 400);
        int x, y;    x = y = 200;
        Color c = Color.black;
        for (Point p : lines)
            g.setColor(c); c = chooseNextColor(c);
            g.drawLine(x, y, p.x, p.y); x = p.x; y = p.y;
    }
}
public synchronized void mouseClicked(MouseEvent e) // Event thread
{ lines.add(new Point(e.getX(), e.getY())); repaint(); }
...
}
```

Interrupts

- An *interrupt* is an event that disrupts the normal flow of control of a program.
- In many systems, interrupts can be totally *asynchronous*, occurring at arbitrary points in a program. The Java developers considered this unwise and arranged that interrupts would occur only at controlled points.
- In Java programs, one thread can interrupt another to inform it that something unusual needs attention:

```
otherThread.interrupt();
```

- But otherThread does not receive the interrupt until it waits: methods wait, sleep (wait for a period of time), join (wait for thread to terminate), and library methods like mailbox deposit and receive.
- An interrupt causes these methods to throw InterruptedException, so typical use is like this:

```
try {  
    msg = inBox.receive();  
} catch (InterruptedException e) { HandleEmergency(); }
```

Remote Mailboxes (A Side Excursion)

- RMI: Remote Method Interface allows one program to refer to objects in another program.
- We use it to allow mailboxes in one program to be received from or be deposited into in another.
- To use this, you define an *interface* to the remote object:

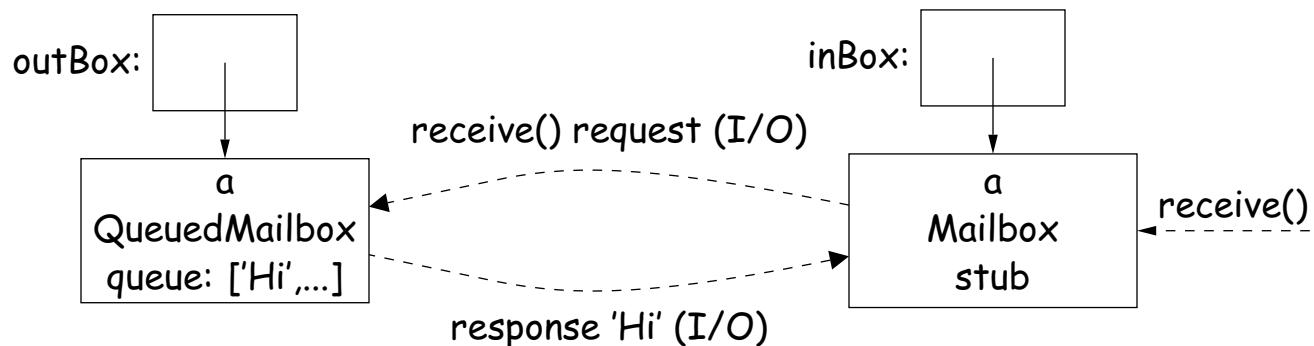
```
import java.rmi.*;  
interface Mailbox extends Remote {  
    void deposit(Object msg)  
        throws InterruptedException, RemoteException;  
    Object receive()  
        throws InterruptedException, RemoteException;  
    ...  
}
```

- On the machine that actually will contain the object, you define

```
class QueuedMailbox ... implements Mailbox {  
    Same implementation as before, roughly  
}
```

Remote Objects Under the Hood

```
// On machine #1:          // On Machine #2:  
Mailbox outBox           Mailbox inBox  
= new QueuedMailbox();    = get outBox from machine #1
```



- Because `Mailbox` is an interface type, you don't see whether you are looking at a mailbox or at a (remote) stub that stands in for it.
- Requests for method calls are relayed by I/O to the machine that has real object.
- Any argument or return type OK if it also implements `Remote` or can be `serialized`—turned into a stream of bytes and back, as can primitive types and `String`.
- Because I/O is involved, expect failures, hence every method can throw `RemoteException` (subtype of `IOException`).

Lecture #38: Storage Management

Scope and Lifetime

- The *scope* of a declaration is portion of program text to which it applies (is *visible*).
 - Need not be a contiguous region.
 - In Java, as in Python, it is static: independent of data.
- The *lifetime* or *extent* of storage is the portion of program execution during which it exists.
 - Always contiguous.
 - Generally dynamic: depends on data
- Classes of extent:
 - *Static*: entire duration of program
 - *Local* or *automatic*: duration of the execution of a call or block—as for local variables or parameters.
 - *Dynamic*: From time of explicit allocation (*new*) to deallocation, if any.

Explicit vs. Automatic Freeing

- Java has no explicit means to free dynamic storage.
- However, when no expression in any thread can possibly be influenced by or change an object, it might as well not exist:

```
IntList wasteful()
{
    IntList c = new IntList(3, new IntList(4, null));
    return c.tail;
    // variable c now deallocated, so no way
    // to get to first cell of list
}
```

- At this point, Java's runtime, like Scheme's, "recycles" the object that c pointed to: *garbage collection*.

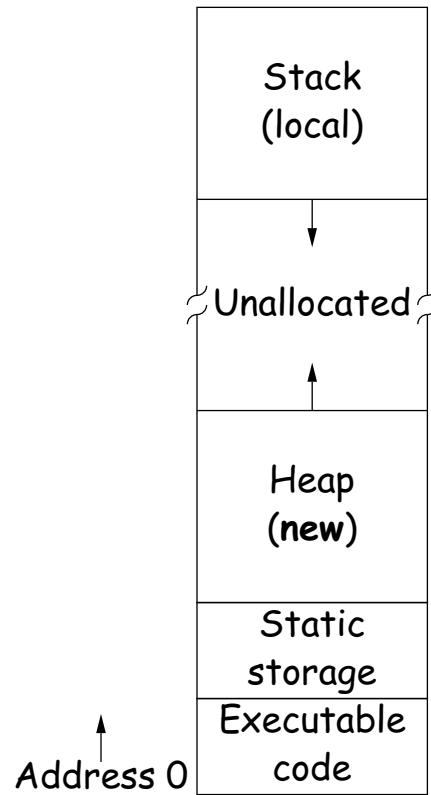
Under the Hood: Allocation

- Java pointers (references) are represented as integer addresses.
- Corresponds to machine's own practice.
- In Java, cannot convert integers \leftrightarrow pointers,
- But crucial parts of Java's runtime are implemented in C, or sometimes machine code, where you can conflate integers and pointers.
- Example of a crude allocator in C:

```
char store[STORAGE_SIZE]; // Allocated array
size_t remainder = STORAGE_SIZE;

/** A pointer to a block of at least N bytes of storage */
void* simpleAlloc(size_t n) { // void*: pointer to anything
    if (n > remainder) ERROR();
    remainder = (remainder - n) & ~0x7; // Make multiple of 8
    return (void*) (store + remainder);
}
```

Example of Storage Layout: Unix



- OS provides a way to turn chunks of unallocated region into heap.
- Happens automatically for stack.

Explicit Deallocating

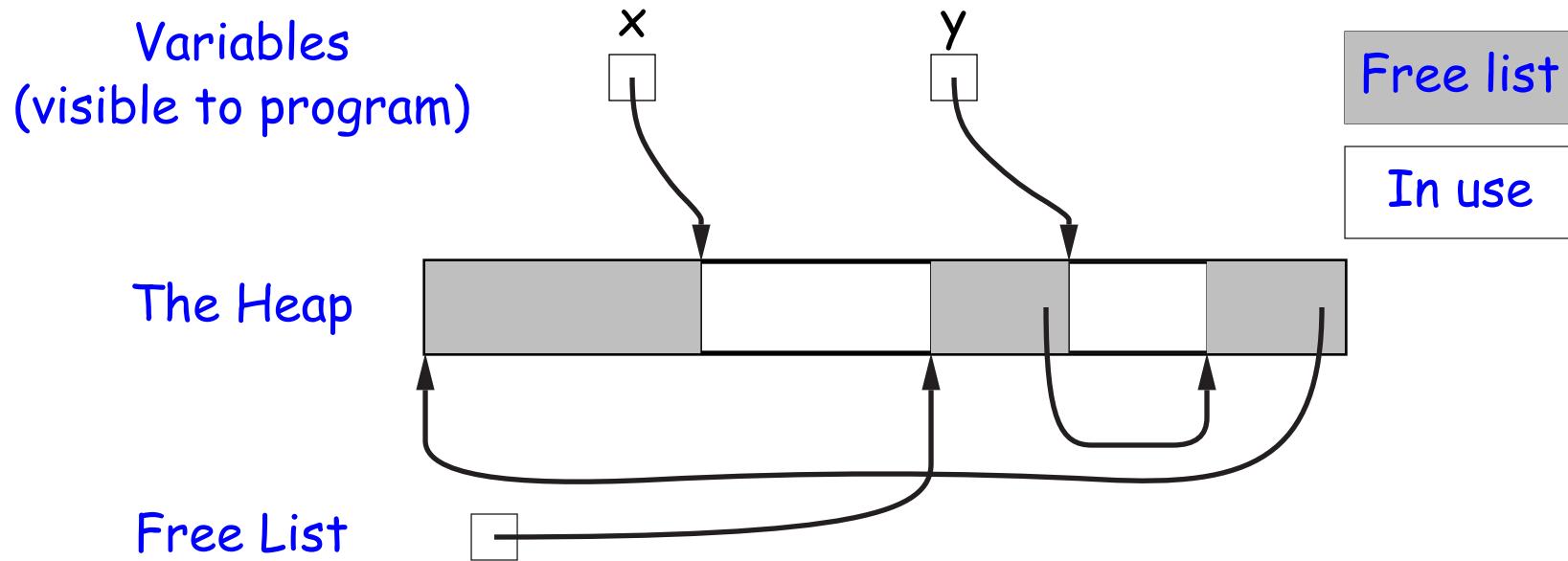
- C/C++ normally require explicit deallocation, because of
 - Lack of run-time information about types and array sizes;
 - Possibility of converting pointers to integers;
 - Lack of run-time information about *unions*:

```
union Various {  
    int Int;  
    char* Pntr;  
    double Double;  
} X; // X is either an int, char*, or double
```

- Java avoids all three problems; automatic collection possible.
- Explicit freeing can be somewhat faster, but rather error-prone:
 - Memory corruption (freeing twice, freeing something that isn't actually a valid pointer.)
 - Memory leaks (failing to ever release something.)

Free Lists

- Explicit allocator grabs chunks of storage from OS to give to applications.
- Or gives recycled storage, when available.
- When storage is freed, it is added to a *free list* data structure to be recycled.
- Used both for explicit freeing and some kinds of automatic storage management.

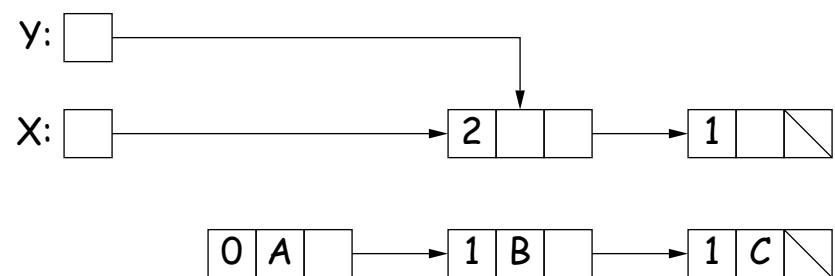
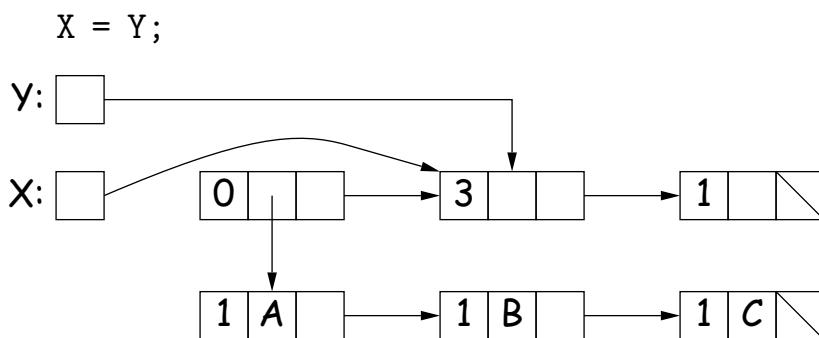
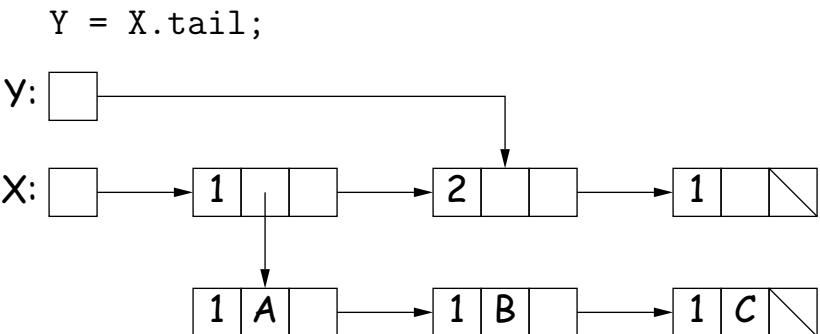
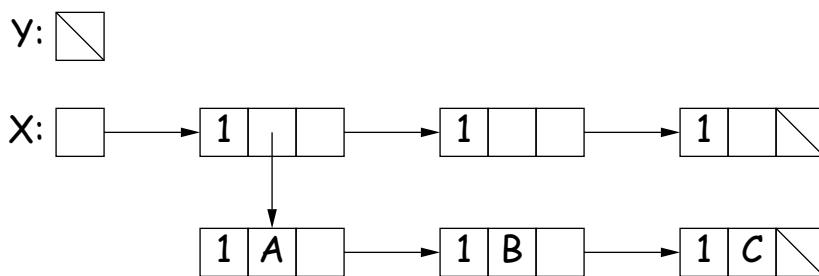


Free List Strategies

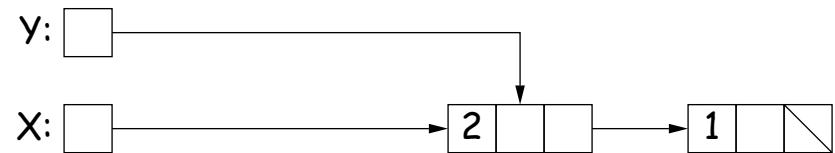
- Memory requests generally come in multiple sizes.
- Not all chunks on the free list are big enough, and one may have to search for a chunk and break it up if too big.
- Various strategies to find a chunk that fits have been used:
 - *Sequential fits:*
 - * Link blocks in LIFO or FIFO order, or sorted by address.
 - * Coalesce adjacent blocks.
 - * Search for *first fit* on list, *best fit* on list, or *next fit* on list after last-chosen chunk.
 - *Segregated fits:* separate free lists for different chunk sizes.
 - *Buddy systems:* A kind of segregated fit where some newly adjacent free blocks of one size are easily detected and combined into bigger chunks.
- Coalescing blocks reduces *fragmentation* of memory into lots of little scattered chunks.

Automatic Garbage Collection: Reference Counting

- Idea: Keep count of number of pointers to each object. Release when count goes to 0.

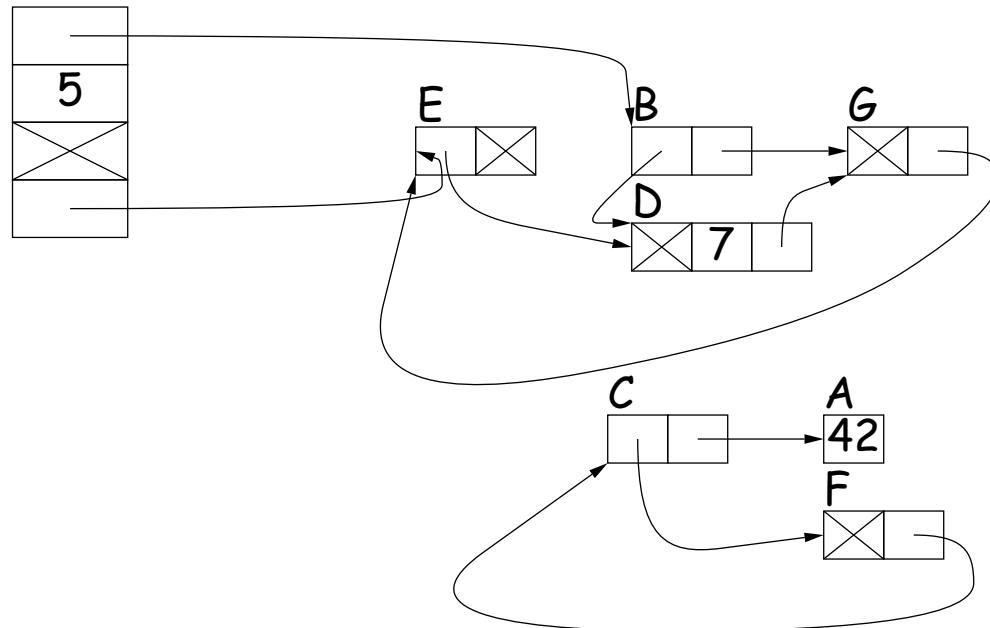


... etc., until:

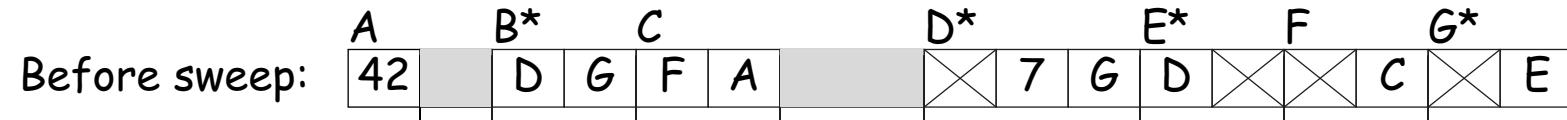


Garbage Collection: Mark and Sweep

Roots (locals + statics)



1. Traverse and **mark** graph of objects.
2. **Sweep** through memory, freeing unmarked objects.



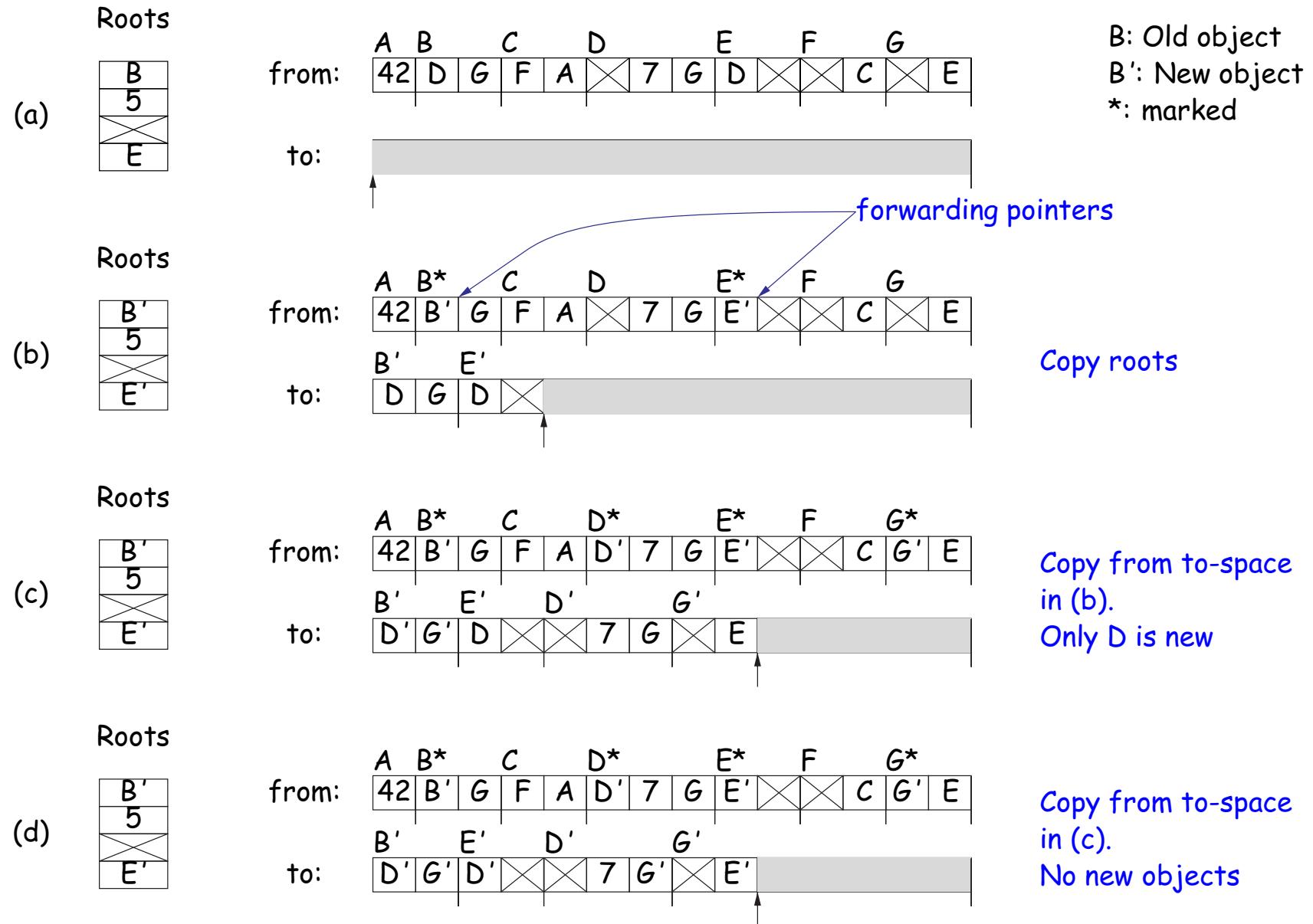
Cost of Mark-and-Sweep

- Mark-and-sweep algorithms don't move any existing objects—pointers stay the same.
- The total amount of work depends on the amount of memory swept—i.e., the total amount of active (non-garbage) storage + amount of garbage. Not necessarily a big hit: the garbage had to be active at one time, and hence there was always some “good” processing in the past for each byte of garbage scanned.

Copying Garbage Collection

- Another approach: *copying garbage collection* takes time proportional to amount of active storage:
 - Traverse the graph of active objects breadth first, *copying* them into a large contiguous area (called "to-space").
 - As you copy each object, mark it and put a *forwarding pointer* into it that points to where you copied it.
 - The next time you have to copy an already marked object, just use its forwarding pointer instead.
 - When done, the space you copied from ("from-space") becomes the next to-space; in effect, all its objects are freed in constant time.

Copying Garbage Collection Illustrated



Most Objects Die Young: Generational Collection

- Most older objects stay active, and need not be collected.
- Would be nice to avoid copying them over and over.
- *Generational garbage collection* schemes have two (or more) from spaces: one for newly created objects (*new space*) and one for "tenured" objects that have survived garbage collection (*old space*).
- A typical garbage collection collects only in new space, ignores pointers from new to old space, and moves objects to old space.
- As roots, uses usual roots plus pointers in old space that have changed (so that they might be pointing to new space).
- When old space full, collect all spaces.
- This approach leads to much smaller *pause times* in interactive systems.

There's Much More

- These are just highlights.
- Lots of work on how to implement these ideas efficiently.
- *Distributed garbage collection*: What if objects scattered over many machines?
- *Real-time collection*: where predictable pause times are important, leads to *incremental* collection, doing a little at a time.

Lecture #39: Compression

Credits: This presentation is largely taken from CS61B lectures by Josh Hug.

Compression and Git

- Git creates a new object in the repository each time a changed file or directory is committed.
- Things can get crowded as a result.
- To save space, it *compresses* each object.
- Every now and then (such as when sending or receiving from another repository), it packs objects together into a single file: a “packfile.”
- Besides just sticking the files together, uses a technique called *delta compression*.

Delta Compression

- Typically, there will be many versions of a file in a Git repository: the latest, and previous edits of it, each in different commits.
- Git doesn't keep track explicitly of which file came from where, since that's hard in general:
 - What if a file is split into two, or two are spliced together?
- But, can guess that files with same name and (roughly) same size in two commits are probably versions of the same file.
- When that happens, store one of them as a pointer to the other, plus a list of changes.

Delta Compression (II)

- So, store two versions

V1

My eyes are fully open to my awful situation.

I shall go at once to Roderick and make him an oration. I shall tell him I've recovered my forgotten moral senses,

V2

My eyes are fully open to my awful situation.

I shall go at once to Roderick and make him an oration.

I shall tell him I've recovered my forgotten moral senses, and don't give twopence halfpenney for any consequences.

as

V1

[Fetch 1st 6 lines from V2]

V2

My eyes are fully open to my awful situation.

I shall go at once to Roderick and make him an oration.

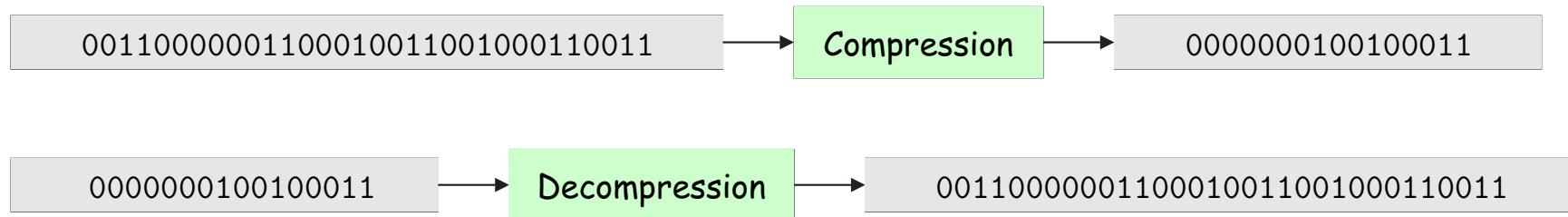
I shall tell him I've recovered my forgotten moral senses, and don't give twopence halfpenney for any consequences.

Two Unix Compression Programs

```
$ gzip -k lect37.pic.in      # The GNU version of ZIP
$ bzip2 -k lect37.pic.in     # Another compression program
$ ls -l lect37.pic*
#
#          Size
#          (bytes)
-rw-r--r-- 1 hilfingr lisp 14794 Apr 25 11:35 lect37.pic.in
-rw-r--r-- 1 hilfingr lisp  5426 Apr 25 11:35 lect37.pic.in.bz2 # Roughly 1/3 size
-rw-r--r-- 1 hilfingr lisp  5529 Apr 25 11:35 lect37.pic.in.gz
$ gzip -k lect37.pdf
$ ls -l lect37.pdf*
-rw-r--r-- 1 hilfingr lisp 79932 Apr 27 11:21 lect37.pdf
-rw-r--r-- 1 hilfingr lisp 66021 Apr 27 11:21 lect37.pdf.gz   # Roughly 83% size
$ gunzip < lect37.pic.in.gz > lect37.pic.in.ungzip # Uncompress
$ diff lect37.pic.in lect37.pic.in.ungzip
$                      # No difference from original (lossless)
$ gzip < lect37.pic.in.gz > lect37.pic.in.gz.gz
$ ls -l lect37.pic*gz
-rw-r--r-- 1 hilfingr lisp 5529 Apr 25 11:35 lect37.pic.in.gz
-rw-r--r-- 1 hilfingr lisp 5552 Apr 27 11:31 lect37.pic.in.gz.gz
$                      # Compressing twice doesn't help.
```

Compression and Decompression

- A *compression algorithm* converts a stream of symbols into another, smaller stream.
- It is called *lossless* if the algorithm is *invertible* (no information lost).
- A common symbol is the bit:



- Here, we simply replaced the 8-bit ASCII bit sequences for digits (where, for example, the single character '0' is encoded as $0x30=0b00110000$) with 4-bit (*binary-coded decimal*).
- Call these 4-bit sequences *codewords*, which we associate with the *symbols* in the original, uncompressed text.
- Can do better than 50% compression with English text.

Example: Morse Code

- Compact, simple to transmit.
- Actually use three symbols:
dih, dah, and pause.
Pauses go between codewords.

A	● —	U	● ● —
B	— ● ● ●	V	● ● ● —
C	— ● — ●	W	● — —
D	— ● ●	X	— ● ● —
E	●	Y	— ● — —
F	● ● — ●	Z	— — — ● ●
G	— — ●		
H	● ● ● ●		
I	● ●	0	— — — — —
J	● — — —	1	● — — — —
K	— ● —	2	● ● — —
L	● — ● ●	3	● ● ● —
M	— —	4	● ● ● ● —
N	— ●	5	● ● ● ● ●
O	— — —	6	— ● ● ● ●
P	● — — ●	7	— — — ● ● ●
Q	— — ● —	8	— — — — ● ●
R	● — ●	9	— — — — — ●
S	● ● ●		
T	—		

Prefix Free Codes

- Morse code needs pauses between codewords to prevent ambiguities.
- Otherwise,

—•●●●——●●●●

could be DEATH, BABE, or BATH.

- The problem is that Morse code allows many codewords to be *prefixes* of other ones, so that it's difficult to know when you have come to the end of one.
- Alternative is to devise *prefix-free codes*, in which no codeword is a prefix of another.
- Then one always knows when a codeword ends.

Prefix-Free Examples

Encoding A

space	1
E	01
T	001
A	0001
O	00001
I	000001
...	

Encoding B

space	111
E	010
T	1000
A	1010
O	1011
I	1100
...	

- For example, "I ATE" is unambiguously
0000011000100101 in Encoding A, or
110011110101000010 in Encoding B.
- What data structures might you use to...
Encode? Decode?

Prefix-Free Examples

Encoding A

space	1
E	01
T	001
A	0001
O	00001
I	000001
...	

Encoding B

space	111
E	010
T	1000
A	1010
O	1011
I	1100
...	

- For example, "I ATE" is unambiguously
0000011000100101 in Encoding A, or
110011110101000010 in Encoding B.
- What data structures might you use to...
Encode? **Ans: HashMap or array** Decode?

Prefix-Free Examples

Encoding A

space	1
E	01
T	001
A	0001
O	00001
I	000001
...	

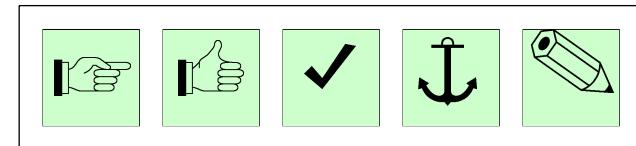
Encoding B

space	111
E	010
T	1000
A	1010
O	1011
I	1100
...	

- For example, "I ATE" is unambiguously
0000011000100101 in Encoding A, or
110011110101000010 in Encoding B.
- What data structures might you use to...
Encode? Ans: HashMap or array Decode? Ans: Trie

Shannon-Fano Coding

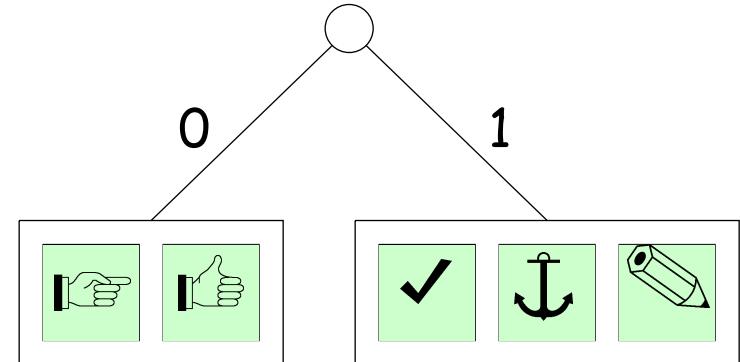
Symbol	Frequency	Encoding
👉	0.35	
👍	0.17	
✓	0.17	
⚓	0.16	
✏️	0.15	



- Count frequencies of all characters in text to be compressed.
- Break grouped characters into two groups of roughly equal frequency.
- Encode left group with leading 0, right group with leading 1.
- Repeat until all groups are of size 1.

Shannon-Fano Coding

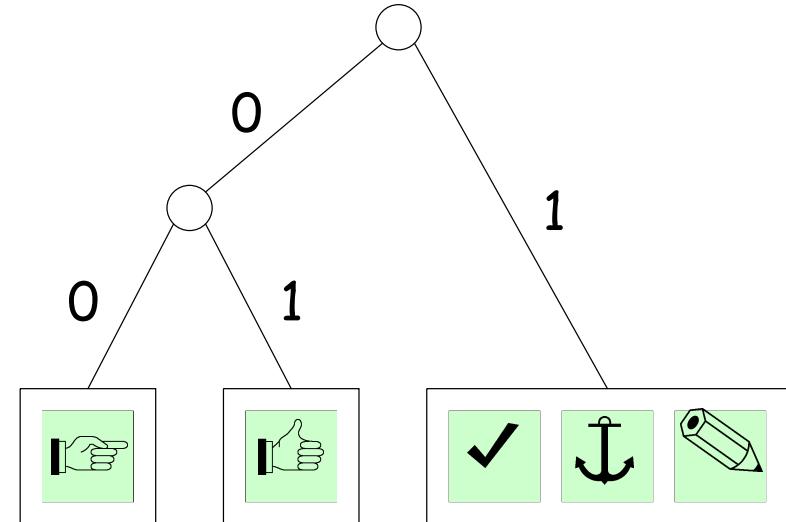
Symbol	Frequency	Encoding
👉	0.35	0...
👍	0.17	0...
✓	0.17	1...
⚓	0.16	1...
-pencil	0.15	1...



- Count frequencies of all characters in text to be compressed.
- Break grouped characters into two groups of roughly equal frequency.
- Encode left group with leading 0, right group with leading 1.
- Repeat until all groups are of size 1.

Shannon-Fano Coding

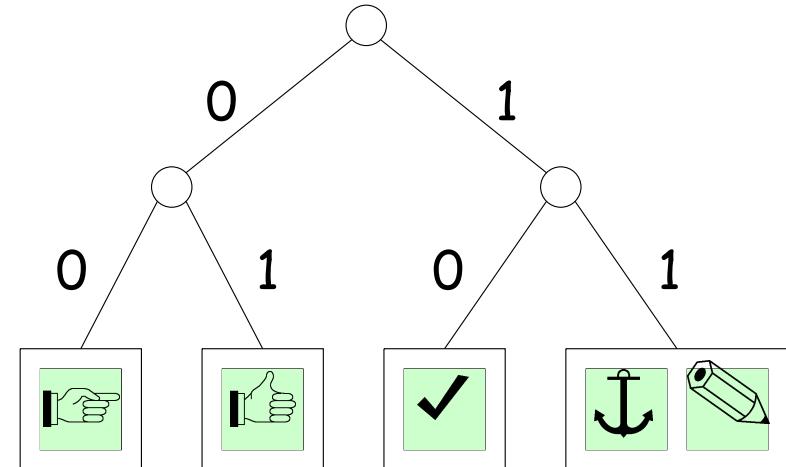
Symbol	Frequency	Encoding
👉	0.35	00
👍	0.17	01
✓	0.17	1...
⚓	0.16	1...
-pencil	0.15	1...



- Count frequencies of all characters in text to be compressed.
- Break grouped characters into two groups of roughly equal frequency.
- Encode left group with leading 0, right group with leading 1.
- Repeat until all groups are of size 1.

Shannon-Fano Coding

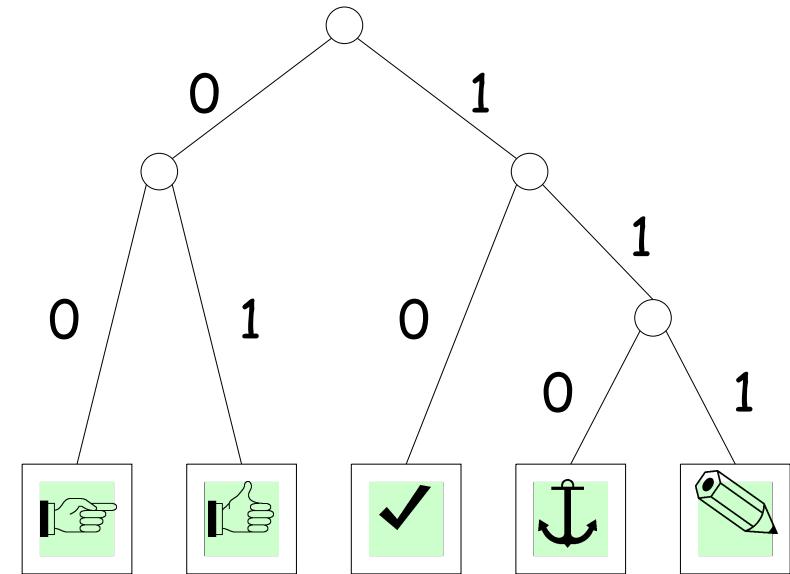
Symbol	Frequency	Encoding
👉	0.35	00
👍	0.17	01
✓	0.17	10
⚓	0.16	11...
-pencil	0.15	11...



- Count frequencies of all characters in text to be compressed.
- Break grouped characters into two groups of roughly equal frequency.
- Encode left group with leading 0, right group with leading 1.
- Repeat until all groups are of size 1.

Shannon-Fano Coding

Symbol	Frequency	Encoding
👉	0.35	00
👍	0.17	01
✓	0.17	10
⚓	0.16	110
✏️	0.15	111

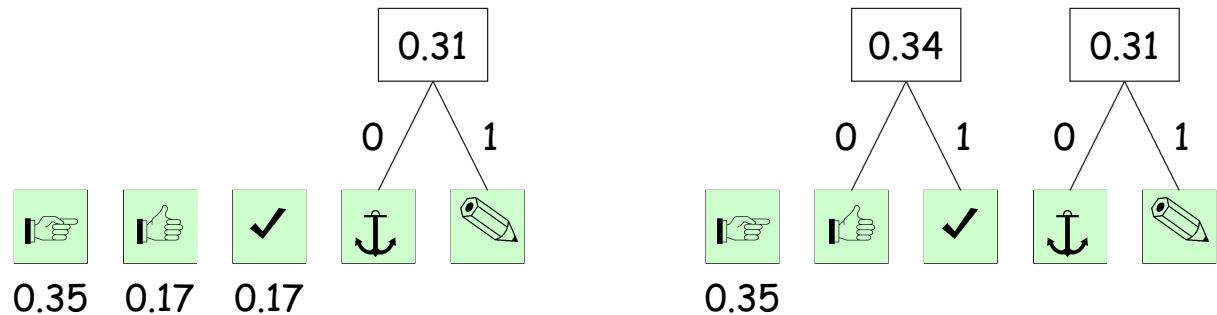


- Count frequencies of all characters in text to be compressed.
- Break grouped characters into two groups of roughly equal frequency.
- Encode left group with leading 0, right group with leading 1.
- Repeat until all groups are of size 1.

Can We Do Better?

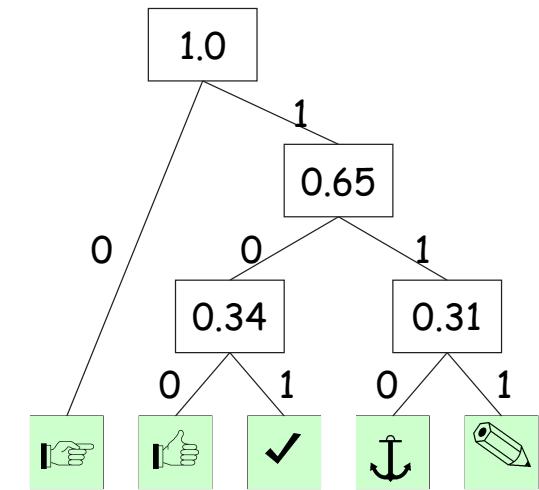
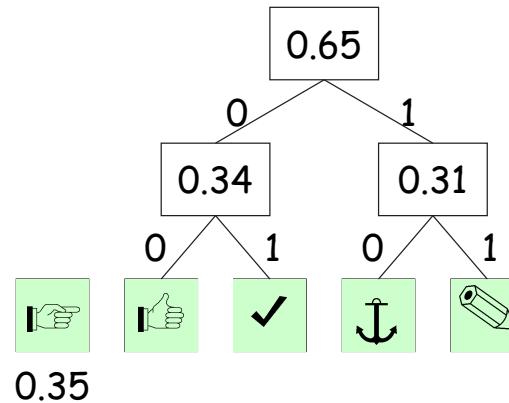
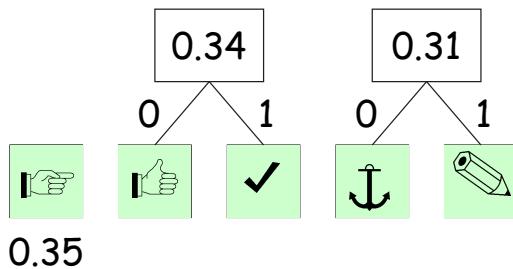
- We'll say an encoding of symbols to codewords that are bitstrings is *optimal* for a particular text if it encodes the text in the fewest bits.
- Shannon-Fano coding is good, but not optimal.
- The optimal solution was found by an MIT graduate student, David Huffman, in a class taught by Fano. The students were given the choice of taking the final or solving this problem (i.e., finding the encoding and a proof of optimality).
- The result is called *Huffman coding*.
- That's right: Fano assigned a problem he hadn't been able to solve. Professors do that occasionally.
- See also this article.

Huffman Coding



- Put each symbol in a node labeled with the symbol's relative frequency (as before).
- Repeat the following until there is just one node:
 - Combine the two nodes with smallest frequencies as children of a new single node whose frequency is the sum of those of the two nodes being combined.
 - Let the edge to the left child be labeled '0' and to the right be labeled '1'.
- The resulting tree shows the encoding for each symbol: concatenate the edge labels on the path from the root to the symbol.

Huffman Coding



- Put each symbol in a node labeled with the symbol's relative frequency (as before).
- Repeat the following until there is just one node:
 - Combine the two nodes with smallest frequencies as children of a new single node whose frequency is the sum of those of the two nodes being combined.
 - Let the edge to the left child be labeled '0' and to the right be labeled '1'.
- The resulting tree shows the encoding for each symbol: concatenate the edge labels on the path from the root to the symbol.

Comparison

Symbol	Frequency	Shannon-Fano	Huffman
	0.35	00	0
	0.17	01	100
	0.17	10	101
	0.16	110	110
	0.15	111	111

For this case, Shannon-Fano coding takes a weighted average of 2.31 bits per symbol, while Huffman coding takes 2.3.

LZW Coding

- So far, we have used systems with one codeword per symbol.
- To get better compression, must encode *multiple* symbols per codeword.
- This will allow us to code strings such as

bbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbb
ababababababababababababababababababab
abcdabcdeabcdefabcdefgabcdefghiabcd

(each 43 characters long) in space that can be less than
43 × weighted average symbol length.

- In LZW coding, we create new codewords as we go along, each corresponding to substrings of the text:
 - Start with a trivial mapping of codewords to single symbols.
 - After outputting a codeword that matches the longest possible prefix, X, of the remaining input, add a new codeword Y that maps to the substring X followed by the next input symbol.

Example of LZW encoding

- Start with a trivial mapping of codewords to single symbols.
- After outputting a codeword that matches the longest possible prefix, X, of the remaining input, add a new codeword Y that maps to the substring X followed by the next input symbol.

Consider the following text as an example:

B="aababcabcdabcdeabcdefabcdefgabcdefgh"

We'll compute $C(B)$, the encoding of B. Our codewords will consist of 8-bit ASCII codes (0x00-0x7f).

LZW Step 0: Initial state

B = aababcabcdabcdeabcdefgabcdefgh

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	

$$C(B) =$$

LZW Step 1

B = aababcabcdabcdeabcdefgabcdefgh

- Best prefix match in the table is 'a', so output 0x61,
- And add aa to the table with a new code.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa

$$C(B) = 0x61$$

LZW Step 2

B = **a**bababcabcdabcdeabcdefgabcdefgh

- Best prefix match in the table for remaining input is still 'a', so output 0x61,
- And add **ab** to the table with a new code.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa
0x81	ab

$$C(B) = 0x6161$$

LZW Step 3

B = aa**b**abcabcdabcdeabcdefgabcdefgh

- Best prefix match in the table for remaining input is 'b', so output 0x62,
- And add **ba** to the table with a new code.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa
0x81	ab
0x82	ba

$$C(B) = 0x616162$$

LZW Step 4

B = **aab**abcabcdabcdeabcdeabcdefabcdefgabcdegh

- Best prefix match in the table for remaining input is now 'ab', so output 0x81 (half as many bits as 'ab').
 - And add $\boxed{ab}c$ to the table with a new code.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa
0x81	ab
0x82	ba
0x83	abc

$$C(B) = 0x61616281$$

LZW Step 5

B = **aabab**cabcdabcdeabcdefabcdeffgabcdeffgh

- Best prefix match in the table for remaining input is now 'c', so output 0x63
 - And add `c a` to the table with a new code.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa
0x81	ab
0x82	ba
0x83	abc
0x84	ca

$$C(B) = 0x6161628163$$

LZW Step 6

B = **aababc**abcdabcdeabcdefgabcdefgh

- Best prefix match in the table for remaining input is now ???, so output ???
- And add ??? to the table with a new code.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa
0x81	ab
0x82	ba
0x83	abc
0x84	ca
0x85	???

$$C(B) = 0x6161628163??$$

LZW Step 6

B = **aababc** **abc** dabcdeabcdefgabcdefgh

- Best prefix match in the table for remaining input is now 'abc', so output 0x83
- And add **abcd** to the table with a new code.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa
0x81	ab
0x82	ba
0x83	abc
0x84	ca
0x85	abcd

$$C(B) = 0x616162816383$$

LZW Step 7

B = **aababcabc** **d** abcdeabcdefcabcdeabcdeabcdeabcde

- Best prefix match in the table for remaining input is now 'd', so output 0x64
- And add 'da' to the table with a new code.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa
0x81	ab
0x82	ba
0x83	abc
0x84	ca
0x85	abcd
0x86	da

$$C(B) = 0x61616281638364$$

LZW Step 7

B = aababcabc d abcdeabcdefcabcdeefgh

- Best prefix match in the table for remaining input is now 'd', so output 0x64
- And add 'da' to the table with a new code.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa
0x81	ab
0x82	ba
0x83	abc
0x84	ca
0x85	abcd
0x86	da

$$C(B) = 0x61616281638364$$

- What's next?
- What is the complete encoding? (When reviewing, try to figure it out before looking at the next slide.)

LZW Final State

$B = \text{aababcabcdabcdeabcdefgabcdefgh}$ (200 bits)

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa
0x81	ab
0x82	ba
0x83	abc
0x84	ca
0x85	abcd
0x86	da

Code	String
0x87	abcde
0x88	ea
0x89	abcdef
0x8a	fa
0x8b	abcdefg
0x8c	ga
0x8d	abcdefgh

$$C(B) = 0x616162816383648565 \\ 876689678b68 \\ (120 \text{ bits})$$

To think about: How might you represent this table to allow easily finding the longest prefix at each step?

Decompression

- Because each different input creates a different table, it would seem that we need to provide the generated table in order to decode a message.
- Interestingly, though, we don't!
- Suppose that, starting with the same initial table we did before, with codes 0x00-0x7f already assigned, we're given

$$C(B) = 0x616162816383$$

and wish to find B .

- We can see it starts with aab. What's next?

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	

Reconstructing the Coding Table (I)

- Idea is to reconstruct the table as we process each codeword in $C(B)$.
- Let $S(X)$ mean “the symbols encoded by codeword X ,” and let Y_k mean character k of string Y .
- For each codeword, X , in $C(B)$, add $S(X)$ to our result.
- Whenever we decoded two consecutive codewords, X_1 and X_2 , **add** a new codeword that maps to $S(X_1) + S(X_2)_0$
- Thus, we recapitulate a step in the compression operation that created $C(B)$ in the first place.
- Since we go from left to right, the table will (almost) always already contain the mapping we need for the next codeword.

LZW Decompression, Step 1

$$C(B) = 0x\boxed{61}616281638364$$

- $S(0x61)$ is 'a' in the table, so add it to B .
- Don't have a previous codeword yet, so don't add anything to the table.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	

$B = a$

LZW Decompression, Step 2

$$C(B) = 0x\text{61}|\boxed{61}6281638364$$

- $S(0x61)$ is 'a' in the table, so add it to B.
- We have two codewords— $S(0x61)$ ='a' twice—so add 'aa' to the table as a new codeword

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa

B = aa

LZW Decompression, Step 3

$$C(B) = 0x\text{6161} \boxed{62} 81638364$$

- $S(0x62)$ is 'b' in the table, so add it to B.
- We have two codewords— $S(0x61)$ ='a' and $S(0x62)$ ='b'—so add 'ab' to the table as a new codeword.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa
0x81	ab

B = aab

LZW Decompression, Step 4

$$C(B) = 0x616162\boxed{81}638364$$

- $S(0x81)$ is 'ab' in the table, so add it to B.
- We have two codewords— $S(0x62)$ ='b' and $S(0x81)$ ='ab'—so add 'ba' to the table as a new codeword.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa
0x81	ab
0x82	ba

B = aabab

LZW Decompression, Step 5

$C(B) = 0x61616281\boxed{63}8364$

- $S(0x63)$ is 'c' in the table, so add it to B.
- We have two codewords— $S(0x81)$ ='ab' and $S(0x63)$ ='c'—so add 'abc' to the table as a new codeword.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa
0x81	ab
0x82	ba
0x83	abc

$B = aababc$

LZW Decompression, Step 6

$$C(B) = 0x6161628163\boxed{83}64$$

- $S(0x83)$ is ??? in the table, so add it to B.
- We have two codewords— $S(???)=???$ and $S(???)=???$ —so add ??? to the table as a new codeword.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa
0x81	ab
0x82	ba
0x83	abc
???	???

B = aababc???

LZW Decompression, Step 6

$$C(B) = 0x6161628163\boxed{83}64$$

- $S(0x83)$ is 'abc' in the table, so add it to B.
- We have two codewords— $S(0x63)$ ='c' and $S(0x83)$ ='abc'—so add 'ca' to the table as a new codeword.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa
0x81	ab
0x82	ba
0x83	abc
0x84	ca

B = aababcabc

LZW Decompression, Step 7

$$C(B) = 0x616162816383 \boxed{64}$$

- $S(0x64)$ is 'd' in the table, so add it to B.
- We have two codewords— $S(0x83)$ ='abc' and $S(0x64)$ ='d'—so add 'abcd' to the table as a new codeword.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	aa
0x81	ab
0x82	ba
0x83	abc
0x84	ca
0x85	abcd

B = aababcabcd

Reconstructing the Coding Table (II)

- In a previous slide, I said "...the table will (almost) always already contain the mapping we need..."
- Unfortunately, there are cases where it doesn't.
- Consider the string $B='cdccdc'$ as an example.
- After we encode it, we end up with

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	cd
0x81	dc
0x82	cdc

$$C(B) = 0x63648082$$

- But decoding causes trouble...

Tricky Decompression, Step 1

$$C(B) = 0x\boxed{63}648082$$

- $S(0x63)$ is 'c' in the table, so add it to B.
- Don't have a previous codeword yet, so don't add anything to the table.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	

$$B = c$$

Tricky Decompression, Step 2

$$C(B) = 0x63\boxed{64}8082$$

- $S(0x64)$ is 'd' in the table, so add it to B.
- We have two codewords— $S(0x63)='c'$ and $S(0x64)='d'$ —so add 'cd' to the table as a new codeword

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	cd

$$B = cd$$

Tricky Decompression, Step 3

$$C(B) = 0x\text{6364} \boxed{80} 82$$

- $S(0x80)$ is 'cd' in the table, so add it to B.
- We have two codewords— $S(0x64)$ ='d' and $S(0x80)$ ='cd'—so add 'dc' to the table as a new codeword

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	cd
0x81	dc

B = cdcd

Tricky Decompression, Step 4

$C(B) = 0x636480[82]$

- Oops! $S(0x82)$ is not yet in the table. What now?

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	cd
0x81	dc
0x82	???

$B = cdcd???$

- Problem is that we could look *ahead* while coding, but can only look *behind* when decoding.
- So must figure out what $0x82$ is going to be by looking back.

Tricky Decompression, Step 4 (Second Try)

$$C(B) = 0x\text{636480} \boxed{82}$$

- $S(0x82)=Z$ (to be figured out).
- Previously decoded $S(0x80) = "cd"$ and now have $S(0x82)=Z$, so will add " cdZ_0 " to the table as $S(0x82)$.
- So Z starts with $S(0x80)$ and therefore Z_0 must be 'c'!
- Thus $S(0x82) = S(0x80)+Z_0 = 'cdc'$.

Code	String
0x61	a
0x62	b
0x63	c
...	...
0x7e	~
0x7f	
0x80	cd
0x81	dc
0x82	cdc

$$B = cdcdcdc$$

LZW Algorithm

- LZW is named for its inventors: Lempel, Ziv, and Welch.
- Was widely used at one time, but because of patent issues became rather unpopular (especially among open-source folks).
- The patents expired in 2003 and 2004.
- Now found in the .gif files, some PDF files, the BSD Unix compress utility and elsewhere.
- There are numerous other (and better) algorithms (such as those in gzip and bzip2).
- The presentation here is considerably simplified.
 - We used fixed-length (8-bit) codewords, but the full algorithm produces variable-length codewords using (!) Huffman coding (compressing the compression).
 - The full algorithm clears the table from time to time to get rid of little-used codewords.

Some Thoughts

- Compressing a compressed text doesn't result in much compression.
- Why must it be impossible to keep compressing a text?
- A program that takes no input and produces an output can be thought of as an encodings of that output.
- Leading to the following question: Given a bitstream, what is the length of the shortest program that can produce it?
- For any specific bitstream, there is a specific answer!
- This is a deep concept, known as Kolmogorov Complexity.

Some Thoughts

- Compressing a compressed text doesn't result in much compression.
- Why must it be impossible to keep compressing a text?
- **Otherwise you'd be able to compress any number of different messages to 1 bit!**
- A program that takes no input and produces an output can be thought of as an encodings of that output.
- Leading to the following question: Given a bitstream, what is the length of the shortest program that can produce it?
- For any specific bitstream, there is a specific answer!
- This is a deep concept, known as Kolmogorov Complexity.

More Thoughts

- It's actually weird that one can compress much at all.
- Consider a 1000-character ASCII text (8000 bits), and suppose we manage to compress it by 50%.
- There are 2^{8000} distinct messages in 8000 bits, but only 2^{4000} possible messages in 4000 bits.
- That is, no matter what one's scheme, one can encode only 2^{-4000} of the possible 8000-bit messsages by 50%! Yet we do it all the time.
- Reason: Our texts have a great deal of redundancy (aka low *information entropy*).
- Texts with high entropy—such as random bits, previously compressed texts, or encrypted texts—are nearly incompressible.

Git

- Git Actually uses a different scheme from LZW for compression: a combination of LZ77 and Huffman coding.
- LZ77 is kind of like delta compression, but within the same text.
- Convert a text such as

One Mississippi, two Mississippi

into something like

One Mississippi, two <11,7>

where the <11,7> is intended to mean "the next 11 characters come from the text that ends 7 characters before this point."

- We add new symbols to the alphabet to represent these (length, distance) inclusions.
- When done, Huffman encode the result.

Lossy Compression

- For some applications, like compressing video and audio streams, it really isn't necessary to be able to reproduce the exact stream.
- We can therefore get more compression by throwing away some information.
- Reason: there is a limit to what human senses respond to.
- For example, we don't hear high frequencies, or see tiny color variations.
- Therefore, formats like JPEG, MP3, or MP4 use *lossy compression* and reconstruct output that is (hopefully) imperceptibly different from the original at large savings in size and bandwidth.
- You can see more of this in EE120 and other courses.

Wrapping Up

- Lossless compression saves space (and bandwidth) by exploiting redundancy in data.
- Huffman and Shannon-Fano coding represent individual symbols of the input with shorter codewords.
- LZW and similar codes represents multiple symbols with shorter codewords.
- Both adapt their codewords to the text being compressed.
- Lossy compression both uses redundancy and exploits the fact that certain consumers of compressed data (like humans) can't really use all the information that could be encoded.

Lecture #40: Course Summary

- Announcement: one EC point for everyone if 75% of you fill out our course survey.
- Programming language: Java
- Program Analysis
- Categories of data structure: Java library structure
- Sequences
- Trees
- Searching
- Sorting
- Pseudo-random numbers
- Graphs
- Pragmatic implementation topics

Programming-Language Topics

- Object-based programming: organizing around data types
- Object-oriented programming:
 - Dynamic vs. static type
 - Inheritance
 - Idea of interface vs. implementation
- Generic programming (the <...> stuff).
- Memory model: containers, pointers, arrays
- Numeric types
- Java syntax and semantics
- Scope and extent
- Standard idioms, patterns:
 - Objects used as functions (e.g., Comparator)
 - Partial implementations (e.g., AbstractList)
 - Iterators
 - Views (e.g., sublists)

Analysis and Algorithmic Techniques

- Asymptotic analysis
- $O(\cdot)$, $o(\cdot)$, $\Omega(\cdot)$, $\Theta(\cdot)$ notations
- Worst case, average case.
- Amortized time
- Memoization and dynamic programming.

Major Categories of Data Structure

- Collection interface and its subtypes
- Map interface and its subtypes
- Generic skeleton implementations of collections, lists, maps (`AbstractList`, etc.)
- Complete concrete collection and map classes in Java library

Sequences

- Linking:
 - Single and double link manipulations
 - Sentinels
- Linking vs. arrays
- Stacks, queues, deques
- Circular buffering
- Trade-offs: costs of basic operations

Trees

- Uses of trees: search, representing hierarchical structures
- Basic operations: insertion, deletion
- Tree traversals
- Representing trees
- Game trees

Searching

- Search trees, range searching
- Multidimensional searches: quad trees.
- Hashing
- Priority queues and heaps
- Balanced trees
 - Rebalancing by rotation (red-black trees)
 - Balance by construction (B-trees)
 - Probabilistic balance (skip lists)
 - Tries
- Search times, trade-offs

Sorting

- Uses of sorting
- Insertion sort
- Selection sorting
- Merge sort
- Heap sort
- Quicksort and selection
- Distribution sort
- Radix sort
- Complexity of various algorithms, when to use them?

Random numbers

- Possible uses
- Idea of a pseudo-random sequence
- Linear congruential and additive generators
- Changing distributions:
 - Changing the range
 - Non-uniform distributions
- Shuffling, random selection

Graph structures

- Definition
- Uses: things represented by graphs
- Graph traversal: the generic traversal template
- Depth-first traversal, breadth-first traversal
- Topological sort
- Shortest paths
- Minimal spanning trees, union-find structures
- Memory management as a graph problem.

Debugging

- What debuggers can do
- How to use to pin down bugs
- Details of some debugger (Eclipse, gjdb, various Windows/Sun products, IntelliJ).
- Unit testing: what it means, how to use it.
- JUnit mechanics.

Version Control

- What's it for?
- Basic concepts behind our particular system:
 - Working copy vs. repository copy
 - Committing changes
 - Updating and merging changes.
 - Tagging

A Case Study

- Presented Git version-control system as an example of a design using several ideas from this course.
- Graph (DAG) and tree structures represented with files as vertices and strings (file names), rather than machine addresses, as pointers.
- Use of hashing to create unique (or very, very likely to be unique) names: *probabilistic data structure*.

Assorted Side Trips

- Compression.
- Parallel processing.
- Storage management and garbage collection.

What's After the Lower Division?

- CSC100/Data C100: Principles & Techniques of Data Science
- CS152: Computer Architecture
- CS160: User Interface Design
- CS161: Computer Security
- CS162: Operating Systems and System Programming
- CS164: Programming Languages and Compilers
- CS168: Intro. to the Internet: Architecture and Protocols
- CS169A, CS169L: Software Engineering

What's After the Lower Division? (II)

- CS170: Efficient Algorithms and Intractable Problems
- CS171: Cryptography
- CS172: Computability and Complexity
- CS174: Combinatorics and Discrete Probability
- CS176: Computational Biology
- CS182: Deep Neural Networks
- CS184: Foundations of Computer Graphics
- CS186: Databases
- CS188: Artificial Intelligence
- CS189: Machine Learning
- CSC191: Quantum Information Science and Technology

What's After the Lower Division? (III)

- CS194: Assorted Special Topics (Examples: Computer Vision and Computational Photography, Parallel Programming)
- CS195: Social Implications of Computer Technology
- Numerous graduate courses: including advanced versions of 152, 160, 161, 170, 184, 186, 189; plus Cryptography, VLSI design and many special topics.
- And, of course, EE courses!
- Various opportunities for participating in research and independent study (199) or directed group studies (198).

And Beyond!

- But EE and CS are just two of over 150 subjects!
- Internships offer more specific skills and exposure to real problems.
- Above all, I think that CS is a creative activity that (to the true artists) ought to be fun!