

Threads

Programs consist of single sequences of instructions. A single sequence is called a *thread* (for "thread of control") in

Programs containing *multiple* threads, which (conceptually) run independently.

To gain program access to threads, Java provides the type `Thread`. Each `Thread` contains information about, and controls,

Access to data from two threads can cause chaos, so Java also constructs for controlled communication, allowing *lock* objects, to *wait* to be notified of events, and to *join* other threads.

Java Mechanics

The actions "walking" and "chewing gum":

```
1 implements Runnable { // Walk and chew gum
  id run()               Thread chomp
    (true) ChewGum(); } = new Thread(new Chewer1());
                          Thread clomp
1 implements Runnable { = new Thread(new Walker1());
  id run()               chomp.start(); clomp.start();
    (true) Walk(); }    }
```

Alternative (uses fact that `Thread` implements `Runnable`):

```
extends Thread {
  run()
  (true) ChewGum(); } Thread chomp = new Chewer2(),
                          clomp = new Walker2();
extends Thread {
  run()
  (true) Walk(); }    chomp.start();
                          clomp.start();
```

Communicating the Hard Way

Getting data is tricky: the faster party must wait for the

slower party. Locks for sending data from thread to thread don't

```
changer {
  value = null;
  receive() {
    r; r = null;
    (r == null)
    = value; }
  = null;
  r;

  deposit(Object data) {
    (value != null) { }
    = data;

    DataExchanger exchanger
    = new DataExchanger();

    -----
    // thread1 sends to thread2 with
    exchanger.deposit("Hello!");

    -----
    // thread2 receives from thread1 with
    msg = (String) exchanger.receive();
  }
}
```

A thread can monopolize the machine while waiting; two threads cutting deposit or receive simultaneously cause chaos.

Lecture #37

Excursions into nitty-gritty stuff: Threads, storage management.

But Why?

On a uniprocessor, only one thread at a time actually runs, but you wait, but this is largely invisible. So why bother with

Threads? Programs always have > 1 thread: besides the main thread, others clean up garbage objects, receive signals, update other stuff.

Programs deal with asynchronous events, it is sometimes convenient to break into subprograms, one for each independent, related event.

How do we insulate one such subprogram from another? We provide a form of modularization.

Organized like this: application is doing some computation while other thread waits for mouse clicks (like 'Stop'), another thread on to updating the screen as needed.

Search engines like search engines may be organized this way, with parallel requests.

Even so, sometimes we *do* have a real multiprocessor.

Avoiding Interference

When a thread has data for another, one must wait for the other

to finish. If two threads use the same data structure, generally only one modifies it at a time; other must wait.

It could happen if two threads simultaneously inserted an element into a linked list at the same point in the list?

They could conceivably execute

```
new ListCell(x, p.next);
```

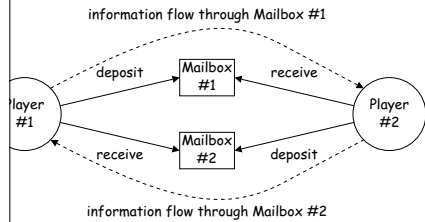
and update the `next` values of `p` and `p.next`; one insertion is lost.

How do we force only one thread at a time to execute a method on an object with either of the following equivalent definitions:

```
} {
synchronized (this) {
  body of f
}
synchronized void f(...) {
  body of f
}
```

Message-Passing Style

Java primitives very error-prone. CS162 goes into alternatives. Higher-level, and allow the following program structure:



Player is a thread that looks like this:

```
nameOver() {
  move()
  tBox.deposit(computeMyMove(lastMove));
}

lastMove = inBox.receive();
```

40:49 2021

CS61B: Lecture #37 8

Coroutines

is a kind of synchronous thread that explicitly hands off to other coroutines so that only one executes at a time, generators. Can get similar effect with threads and

recursive inorder tree iterator:

```
Cor extends Thread {
  Mailbox r;
  (Tree T, Mailbox r) {
    t = T; this.dest = r;
  }
  void treeProcessor(Tree T) {
    Mailbox m = new QueuedMailbox();
    new TreeIterator(T, m).start();
    while (true) {
      Object x = m.receive();
      if (x is end marker)
        break;
      do something with x;
    }
  }
  (Tree t) {
    null return;
    (t.left);
    (t.label);
    (t.right);
  }
}
```

40:49 2021

CS61B: Lecture #37 10

Highlights of a GUI Component

```
that draws multi-colored lines indicated by mouse. */
extends JComponent implements MouseListener {
  <Point> lines = new ArrayList<Point>();

  // Main thread calls this to create one
  setSize(new Dimension(400, 400));
  setListener(this);

  @Override
  void paintComponent(Graphics g) { // Paint thread
    g.setColor(Color.white); g.fillRect(0, 0, 400, 400);
    x = y = 200;
    Color.black;
    p : lines
    g.setColor(c); c = chooseNextColor(c);
    g.fillRect(x, y, p.x, p.y); x = p.x; y = p.y;
  }

  @Override
  void mouseClicked(MouseEvent e) // Event thread
  {
    lines.add(new Point(e.getX(), e.getY())); repaint(); }
}
```

40:49 2021

CS61B: Lecture #37 12

Primitive Java Facilities

wait() method makes the current thread wait (not using processor) until notified by notifyAll(), unlocking the Object while it waits.

java.util.concurrent.locks.Lock has something like this (simplified):

```
interface Mailbox {
  void deposit(Object msg) throws InterruptedException;
  Object receive() throws InterruptedException;
}
```

```
Mailbox implements Mailbox {
  List<Object> queue = new LinkedList<Object>();
}
```

```
public synchronized void deposit(Object msg) {
  queue.add(msg);
  notifyAll(); // Wake any waiting receivers
}
```

```
public synchronized Object receive() throws InterruptedException {
  while (queue.isEmpty()) wait();
  return queue.remove(0);
}
```

40:49 2021

CS61B: Lecture #37 7

More Concurrency

wait() can be done other ways, but mechanism is very similar

you want to think during opponent's move:

```
nameOver() {
  move()
  tBox.deposit(computeMyMove(lastMove));
}
```

```
{
  thinkAheadALittle();
  lastMove = inBox.receiveIfPossible();
  while (lastMove == null);
}
```

receiveIfPossible() (written receive() in our actual package) doesn't return null if no message yet, perhaps like this:

```
public synchronized Object receiveIfPossible()
  throws InterruptedException {
  while (queue.isEmpty())
    return null;
  return queue.remove(0);
}
```

40:49 2021

CS61B: Lecture #37 9

Use In GUIs

Swing library uses a special thread that does nothing but updates the GUI like mouse clicks, pressed keys, mouse movement,

you designate an object of your choice as a *listener*; which Java's event thread calls a method of that object whenever events occur.

your program can do work while the UI continues to update buttons, menus, etc.

the special thread does all the drawing. You don't have to wait for this to take place; just ask that the thread wake up when you change something.

40:49 2021

CS61B: Lecture #37 11

Note Mailboxes (A Side Excursion)

The Method Interface allows one program to refer to another program.

allow mailboxes in one program to be received from or into in another.

you define an *interface* to the remote object:

```
import java.rmi.*;
Mailbox extends Remote {
    deposit(Object msg)
    throws InterruptedException, RemoteException;
    receive()
    throws InterruptedException, RemoteException;
```

the line that actually will contain the object, you define

```
class Mailbox ... implements Mailbox {
// implementation as before, roughly
```

40:49 2021

CS61B: Lecture #37 14

New Topic: Storage Management

Explicit vs. Automatic Freeing

explicit means to free dynamic storage. when no expression in any thread can possibly be influenced by an object, it might as well not exist:

```
list.add(steptail());
```

```
list.add(new IntList(3, new IntList(4, null)));
list.tail();
// variable c now deallocated, so no way to get to first cell of list
```

But, Java's runtime, like Scheme's, "recycles" the object and so does not do: *garbage collection*.

40:49 2021

CS61B: Lecture #37 18

Interrupts

An interrupt is an event that disrupts the normal flow of control of a program.

In Java, interrupts can be totally *asynchronous*, occurring at arbitrary points in a program. The Java developers considered it important to arrange that interrupts would occur only at controlled points.

In programs, one thread can interrupt another to inform it that it has unusual needs attention:

```
thread.interrupt();
// thread does not receive the interrupt until it waits: methods like wait() (wait for a period of time), join() (wait for thread to finish), and library methods like mailbox.deposit() and mailbox.receive().
```

Calling interrupt() causes these methods to throw InterruptedException, and the response is like this:

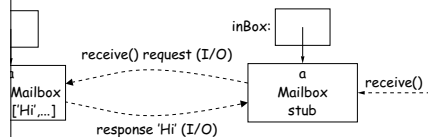
```
try { inBox.receive();
} catch (InterruptedException e) { HandleEmergency(); }
```

40:49 2021

CS61B: Lecture #37 13

Remote Objects Under the Hood

```
// #1: // On Machine #2:
Mailbox inBox;
inBox = remoteMailbox(); // = get outBox from machine #1
```



The Mailbox interface is an interface type, you don't see whether you are dealing with a Mailbox or at a (remote) stub that stands in for it.

Method calls are relayed by I/O to the machine that implements the interface.

The Mailbox interface is OK if it also implements Remote or Serializable—turned into a stream of bytes and back, as can be done with Serializable and String.

When RemoteException is involved, expect failures, hence every method can throw RemoteException (subtype of IOException).

40:49 2021

CS61B: Lecture #37 15

Scope and Lifetime

The scope of a declaration is the portion of program text to which it applies (e.g., *local* or *global*).

The lifetime of a variable is the portion of program execution in which it is active. It is a contiguous region. In Python, it is static: independent of data. In Java, the lifetime of a variable is the portion of program execution in which it exists.

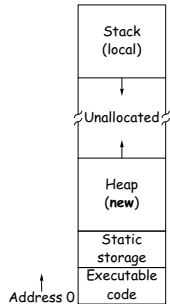
Static: independent of data. Dynamic: depends on data. Lifetime: duration of program execution. *Automatic*: duration of the execution of a call or block—also applies to variables or parameters.

From time of explicit allocation (new) to deallocation,

40:49 2021

CS61B: Lecture #37 17

Example of Storage Layout: Unix



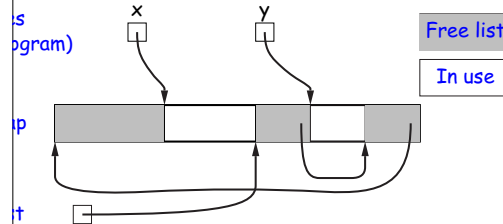
One way to turn chunks of unallocated region into heap. Automatically for stack.

Free Lists

Allocator grabs chunks of storage from OS to give to applications. Recycled storage, when available.

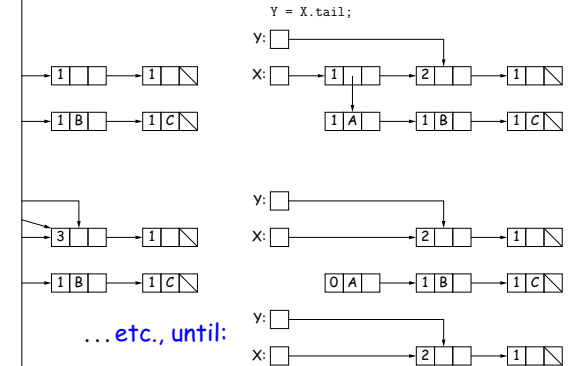
When memory is freed, it is added to a *free list* data structure to

allow for explicit freeing and some kinds of automatic storage



Garbage Collection: Reference Counting

Keep count of number of pointers to each object. Release objects when count goes to 0.



Under the Hood: Allocation

References are represented as integer addresses.

Due to machine's own practice.

Cannot convert integers ↔ pointers,

parts of Java's runtime are implemented in C, or sometimes C++, where you can conflate integers and pointers.

Simplest crude allocator in C:

```
[STORAGE_SIZE]; // Allocated array
pointer = STORAGE_SIZE;
```

```
pointer to a block of at least N bytes of storage */
void* malloc(size_t n) { // void*: pointer to anything
    if (n > remainder) ERROR();
    pointer = (remainder - n) & ~0x7; // Make multiple of 8
    return (void*)(store + remainder);
}
```

Explicit Deallocating

Applications require explicit deallocation, because of

lack of run-time information about types and array sizes;

lack of converting pointers to integers;

lack of run-time information about unions:

```
Various {
    int Int;
    char* Pntr;
    double Double;
    // X is either an int, char*, or double
}
```

Addresses of all three problems; automatic collection possible.

Explicit deallocation can be somewhat faster, but rather error-prone:

Double-free (freeing twice, freeing something that isn't yours), double-free (freeing twice, freeing something that isn't yours), double-free (freeing twice, freeing something that isn't yours).

Memory leaks (failing to ever release something.)

Free List Strategies

Free lists generally come in multiple sizes.

Free lists on the free list are big enough, and one may have to chunk and break it up if too big.

Strategies to find a chunk that fits have been used:

Best fit:

Checks in LIFO or FIFO order, or sorted by address.

Free adjacent blocks.

Searches for *first fit* on list, *best fit* on list, or *next fit* on list to find next-chosen chunk.

Best fit: separate free lists for different chunk sizes.

Best fits: A kind of segregated fit where some newly adjacent blocks of one size are easily detected and combined into larger blocks.

Free lists reduces *fragmentation* of memory into lots of small free chunks.

Cost of Mark-and-Sweep

Sweep algorithms don't move any existing objects—pointers are free.

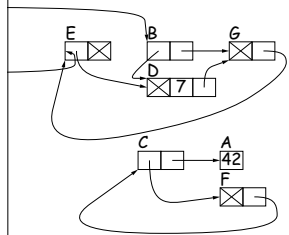
Amount of work depends on the amount of memory swept—i.e., amount of active (non-garbage) storage + amount of garbage. Only a big hit: the garbage had to be active at one time, there was always some "good" processing in the past for garbage scanned.

Copying Garbage Collection Illustrated

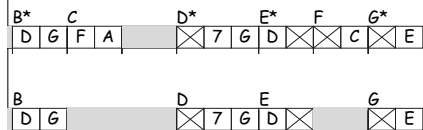


Garbage Collection: Mark and Sweep

steps



1. Traverse and **mark** graph of objects.
2. **Sweep** through memory, freeing unmarked objects.



Copying Garbage Collection

Approach: **copying garbage collection** takes time proportional to amount of active storage:

Traverse the graph of active objects breadth first, **copying** them to a contiguous area (called "to-space").

Copy each object, mark it and put a **forwarding pointer** that points to where you copied it.

Next time you have to copy an already marked object, just use the forwarding pointer instead.

Finally, the space you copied from ("from-space") becomes the new to-space; in effect, all its objects are freed in constant time.

There's Much More

Next highlights.

Focus on how to implement these ideas efficiently.

Scattered garbage collection: What if objects scattered over many memory spaces?

Incremental garbage collection: where predictable pause times are important, doing a little at a time.

Objects Die Young: Generational Collection

Young objects stay active, and need not be collected.

Goal: to avoid copying them over and over.

Generational garbage collection schemes have two (or more) spaces: one for newly created objects (**new space**) and one for objects that have survived garbage collection (**old space**).

Generational garbage collection collects only in new space, ignores pointers in old space, and moves objects to old space.

Generational garbage collection has usual roots plus pointers in old space that have changed (and might be pointing to new space).

When new space full, collect all spaces.

Generational garbage collection leads to much smaller **pause times** in interactive systems.