



Intro to Parallelism

CSE 332 25Sp
Lecture 19

Announcements

	Monday	Tuesday	Wed	Thursday	Friday
This Week	Ex 7 (DFS, coding) due <u>Ex 9</u> (reductions, gs) out		TODAY	Bring <u>laptop!</u>	Ex 8 (Dijkstra, gs) due Ex <u>10,11</u> (parallel prog) out
Next Week	<u>Ex 9</u> (reductions, gs) due				Ex 10 (parallel, prog) due

Bring Laptops to section tomorrow!

{ Optional readings (Grossman) covers next few weeks of parallelism and concurrency }

Parallelism & Concurrency

All of your programs have made the same assumption

One thing happens at a time

Usually called "sequential programming"

Over the next two weeks we'll remove this assumption

- Write programs that divide work between multiple threads and **synchronize** their behavior
- Design algorithms to provide a speedup
 - More throughput: work done per unit time
- Data structures decide how to allow concurrent access to data among all the threads.

Why are we doing this?

Parallelism is where computation is heading.

From 1980-2005 (ish) desktop computers got twice as fast every 18 months or so.

- Moore's Law. Not an immutable law of nature. Business decision.
- How? Keep making everything smaller

Code not running fast enough? It'll be four times as fast if you just buy a new computer.

Why are we doing this?

End of Moore's Law

We're at the limit of our ability to shrink processors.

- Transistors are really small (much smaller and quantum mechanics kicks in)
- and get really hot.

Computer Architects are working very hard to still speed up processors just a little bit more.

- Take an architecture class to get a taste.

But to really achieve a speedup, the solution has been more processors.

Why are we doing this?

Parallelism is where the world is heading.

Our computers are still getting faster by adding more processors

- Rather than just making each new one twice as fast.

If we want to solve new, bigger problems, we're going to need to take advantage of more than one processor.

We won't forget about sequential/single processor programming.

- It will still be simpler and good enough most of the time.

But understanding parallelism is more important than ever.

Parallelism vs. Concurrency

Parallelism: Use extra resources (i.e. processors) to solve your problem faster

Concurrency: Correctly and efficiently sharing a single resource among multiple threads.

Terms aren't completely standard.

They overlap somewhat.

Analogies

Cooking:

Parallelism (do one job faster with more power):

I have hundreds of potatoes to slice.

Get 20 extra cooks (and knives)

Hand them all a bunch of potatoes

Concurrency (manage shared resources):

10 cooks are trying to share 4 burners

And one oven

Examples

Parallelism:

I want to sum up all the elements in an array

Divide the array in 4, sum up each piece in a different thread

Add together the threads' answers for the final answer

Concurrency:

Two users are trying to add an entry to a hash table at the same time.

What if the hashes collide? What if they're the same key and different values?

Sharing Memory with Threads

Our parallelism model will be shared memory with threads.

- There are other models (see Grossman), we won't use them.

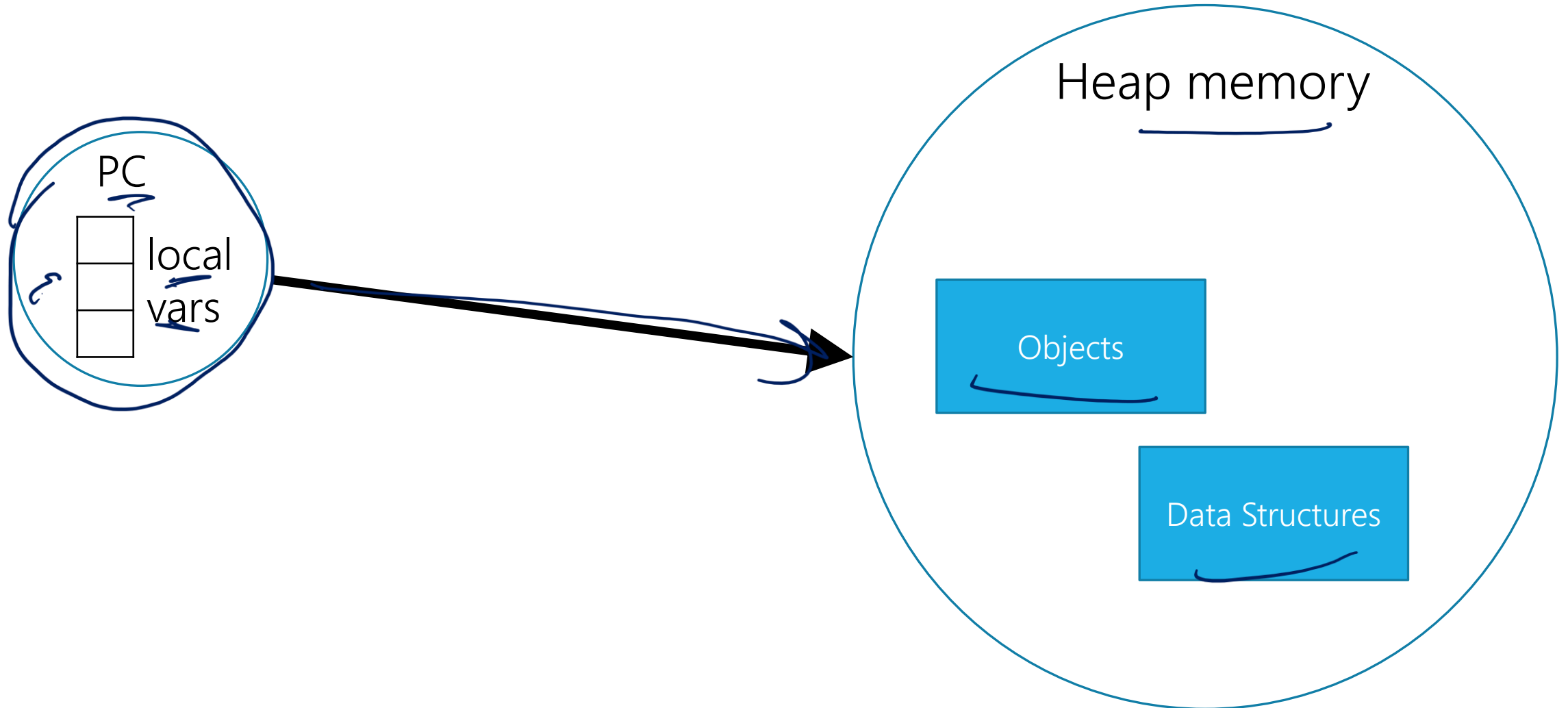
Sequential Story:

- One program counter
- One call stack
- new Objects go in the heap

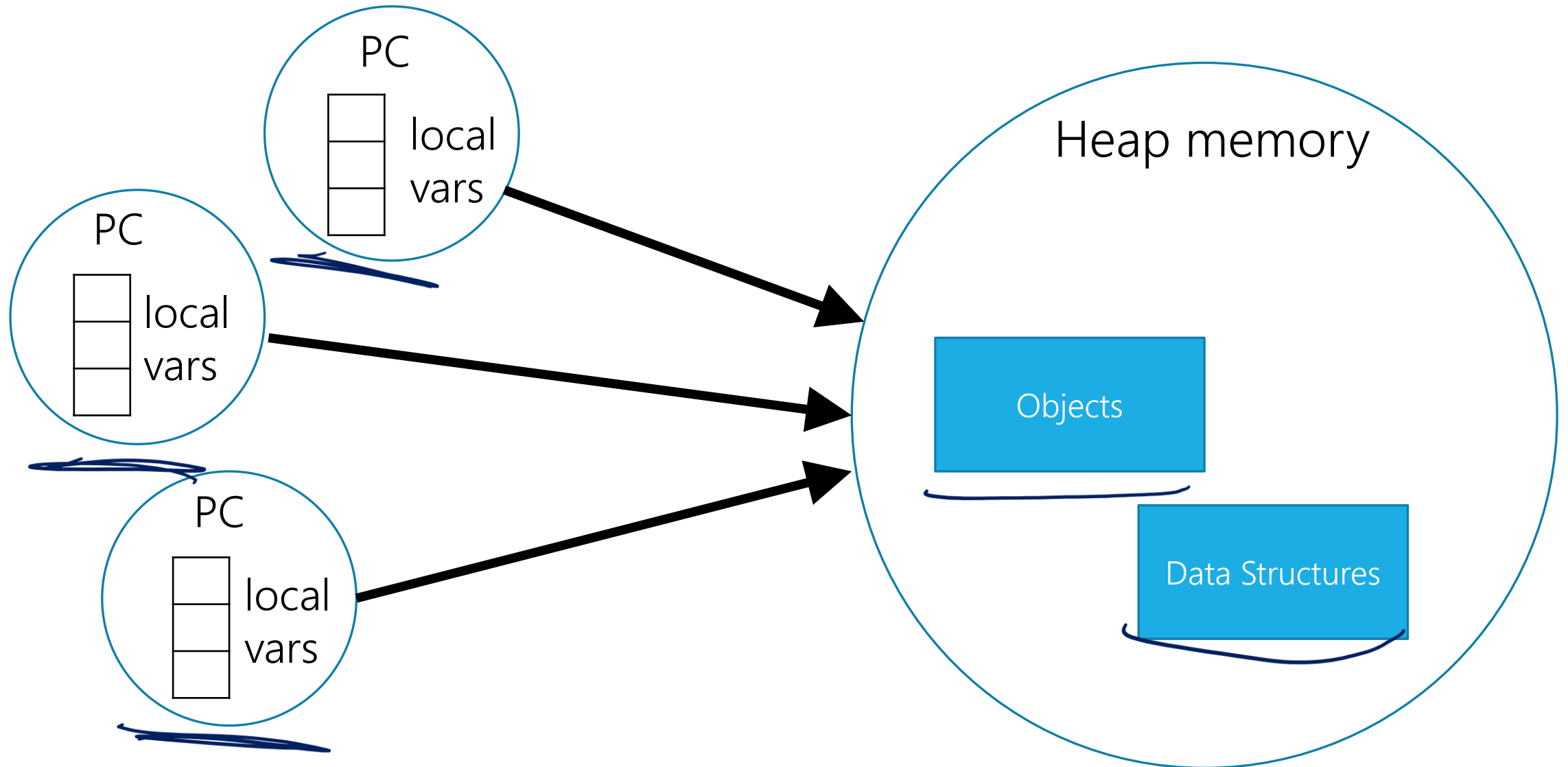
Parallel Story

- Set of threads. Each has its own program counter and its own stack
- Threads will (implicitly) share objects and static fields
- Threads communicate by altering memory.

Sequential Code



Parallel Code



We need new primitives

To write parallel programs we need a library with:

Ways to create and run multiple things at once

- i.e. threads

Ways for threads to share memory

- Usually just having the same references

Ways for threads to coordinate

- This week: A way for threads to wait for others to finish

- Next week: prevent others from accessing memory until we're done

For Today

We'll only write pseudocode (we'll introduce the library soon)

Parallelism requires a different mode of thinking

Just going to practice that on an example problem

A Simple Problem

Goal: Given an array, sum up all the elements.

First idea: Start up 4 threads. Each sums $\frac{1}{4}$ of the array.
Then add together those answers.

ParallelSum: Take 1 (not correct)

```
Class SumThread extends SomeThreadObject{  
    int lo; int hi; int[] arr;  
    int ans = 0; //result  
    SumThread(int[] a, int l, int h){  
        lo = l; hi=h; arr=a;  
    }  
    void run(){  
        for(i=lo; i<hi; i++)  
            ans+=arr[i];  
    }  
}
```


ParallelSum: Take 1

There are major bugs with this code.
Find some of them!

```
int sum(int[] arr) {  
    int len = arr.length;  
    int ans = 0;  
    SumThread[] ts = new SumThread[4];  
    for(int i=0; i<4; i++)  
        ts[i] = new SumThread(arr, i*len/4, (i+1)*len/4);  
    {  
        for(int i=0; i<4; i++)  
            ans += ts[i].ans;  
    }  
    return ans;  
}
```

Bugs

We made some Thread objects...

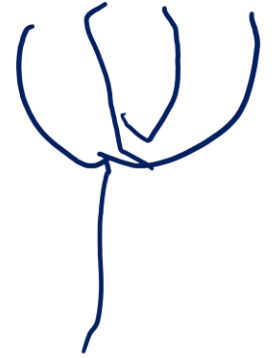
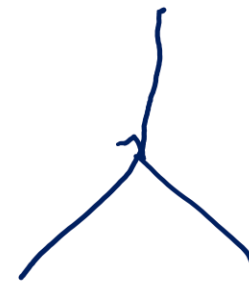
- but we never actually started them. They're just sitting there.
- Be careful what method you call!
- Libraries will have different methods for

→ "look at this thread object, run the code IN YOURSELF not in that thread."

→ "look at this object, tell THAT THREAD to run its code."

ParallelSum: Take 2

```
int sum(int[] arr) {  
    int len = arr.length;  
    int ans = 0;  
    SumThread[] ts = new SumThread[4];  
    for(int i=0; i<4; i++)  
        ts[i] = new SumThread(arr, i*len/4, (i+1)*len/4);  
        ts[i].fork();  
    for(int i=0; i<4; i++)  
        ans += ts[i].ans;  
    return ans;  
}
```



Bugs

We made some Thread objects...

- but we never actually started them. They're just sitting there.
- Be careful what method you call!
- Libraries will have different methods for
 - "look at this thread object, run the code IN YOURSELF not in that thread."
 - "look at this object, tell THAT THREAD to run its code."

The current thread is still running.

Will each thread update its ans field in time?

Need to tell original thread to WAIT for its children to finish.

ParallelSum: Take 2

```
int sum(int[] arr) {  
    int len = arr.length;  
    int ans = 0;  
    SumThread[] ts = new SumThread[4];  
    for(int i=0; i<4; i++)  
        ts[i] = new SumThread(arr, i*len/4, (i+1)*len/4);  
    ts[i].fork()  
    for(int i=0; i<4; i++)  
        ts[i].join()  
    ans += ts[i].ans;  
    return ans;  
}
```

Join

Parallelism libraries will define methods you can't implement on your own.

- E.g. whatever method starts a new thread isn't something you can do yourself.

Join is our first taste of coordinating computation

- Calling thread "blocks" (just sits there doing nothing) until receiver returns
- Avoids race condition in our original code.

This style of programming is called "fork/join"

- Java note: `join` can throw exceptions. May not compile unless you catch a `java.lang.InterruptedException`
- A simple try-catch block should be fine for simple code.

(Almost) No Shared Memory!

Fork/join programs like these don't really share memory

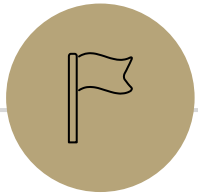
We divided up the array – no one tried to access the same locations.

Lo, hi, arr fields weren't shared. Each helper thread had those values written by the main thread.

Main thread gets data back – doesn't let the helper threads alter any shared data themselves.

To avoid race conditions, we'll use `join`

Next week, we'll see other ways to synchronize.



Optimizations

Optimizing: Number of Threads

The last version of ParallelSum will work.

I.e. it will always get the right answer

And it will use 4 threads.

But...

What if we get a new computer with 6 processors?

- We'll have to rewrite our code

What if the OS decides "no, you only get 2 processors right now."

What if our threads take wildly different amounts of time?

Optimizing: Number of Threads

The counter-intuitive solution:

→ Even more parallelism!

Divide the work into more smaller pieces.

If you get more processors, you take advantage of all of them.

If one thread finishes super fast, throw the next thread to that processor.

"Load Imbalance"

Engineering Question:

- Let's say we change our `ParallelSum` code so each thread adds 10 elements.
- Is that a good idea? What's the running time of the code going to be?

Thread Creation

If we create $n/10$ threads, each summing 10 elements...

Creating $n/10$ threads one-right-after-the-other takes $\Theta(n)$ time.

(Same with joining the threads together at the end).

This is a linear time algorithm now. Can we do better?

Divide and Conquer: Parallelism

What if we want a bunch of threads, but don't want to spend a bunch of time making threads?

Parallelize thread creation too!

Divide and Conquer SumThread

```
Class SumThread extends SomeThreadObject{
    //constructor, fields unchanged.
    void run(){
        if(hi-lo == 1)
            ans = arr[lo]
        else{
            SumThread left = new SumThread(arr, lo, (hi+lo)/2);
            SumThread right = new SumThread(arr, (hi+lo)/2, hi);
            left.start(); right.start();
            left.join(); right.join();
            ans = left.ans + right.ans;
        }
    }
}
```

Divide And Conquer SumThread

```
int sum(int[] arr) {  
    SumThread t = new SumThread(arr, 0, arr.length);  
    t.run(); //this call isn't making a new thread  
    return t.ans;  
}
```

Divide And Conquer Optimization

Imagine calling our current algorithm on an array of size 4.

How many threads does it make

6

It shouldn't take that many threads to add a few numbers.

And every thread introduces A LOT of overhead.

We'll want to **cut-off** the parallelism when the threads cause too much overhead.

Similar optimizations can be used for (sequential) merge and quick sort

Cut-offs

Are we really saving that much?

Suppose we're summing an array of size 2^{30}

And we set a cut-off of size-100

-i.e. subarrays of size 100 are summed without making any new threads.

What fraction of the threads have we just eliminated?

99% !!!! (for fun you should check the math)

One more optimization

A small tweak to our code will eliminate half of our threads

Old

```
SumThread left = new SumThread(arr, lo, (hi+lo)/2);  
SumThread right = new SumThread(arr, (hi+lo)/2, hi);  
left.start(); right.start();  
left.join(); right.join();
```

Better

```
SumThread left = new SumThread(arr, lo, (hi+lo)/2);  
SumThread right = new SumThread(arr, (hi+lo)/2, hi);  
left.start();  
right.run();  
left.join();
```

Order of these lines
matters!

Wrap Up

None of our optimizations will make a difference in the $O()$ analysis

- Which we'll see next time

But they will make a difference in practice.

Next Time:

Using a real library

Analyzing parallel programs.

ForkJoin Framework

Method	Use
<code>compute</code>	Thread objects override (void) method <code>run</code> When <code>fork()</code> is called, <code>run()</code> method is executed A bit like <code>main(String[] args)</code> ---default starting point
<code>fork</code>	Starts a new thread executing that object's <code>run</code> method
<code>join</code>	Calling <code>otherThread.join()</code> pauses <u>this</u> thread until <code>otherThread</code> has completed its <code>run</code> method.
<code>RecursiveTask<V></code>	Class which we extend to make threads