

Practical Concurrent and Parallel Programming 5

Riko Jacob
IT University of Copenhagen

Friday 2018-09-28

Reminder: Locking can be tricky

What is wrong with TestStaticCounter.java?

```
static Long count = 0L
```

```
..
```

```
    synchronized(count) { count++; }
```

```
..
```

Plan for today

- Tasks and the Java executor framework
 - Executors, Runnables, Callables, Futures
 - The states of a task
 - Task creation overhead
 - Using tasks to count prime numbers
 - Java versus the .NET Task Parallel Library
 - Producer-consumer pipelines
 - Bounded queues, thread wait() and notify()
 - The states of a thread
 - Java 8 stream implementation
- Based on slides by Peter Sestoft

Prefer executors and tasks to threads

- We have used *threads* to parallelize work
 - But creating many threads takes time and memory
- Better divide work into (many small) *tasks*
 - Then submit the tasks to an executor
 - This uses a pool of (few) threads to run the tasks
- Goetz chapters 6, 8 and Bloch item 68

should generally refrain from working directly with threads. The key abstraction is no longer `Thread`, which served as both the unit of work and the mechanism for executing it. Now the unit of work and mechanism are separate. The key abstraction is the unit of work, which is called a *task*. There are two kinds of tasks: `Runnable` and its close cousin, `Callable` (which is like `Runnable`, except that it returns a value). The general mechanism for executing tasks is the *executor ser-*

Bloch item 68

Executors and tasks

- A *task* is just a `Runnable` or `Callable<T>`
- Submitting it to an *executor* gives a *Future*

```
Future<?> fut
= executor.submit(new Runnable() { public void run() {
    System.out.println("Task ran!");
}});
```

TestTaskSubmit.java

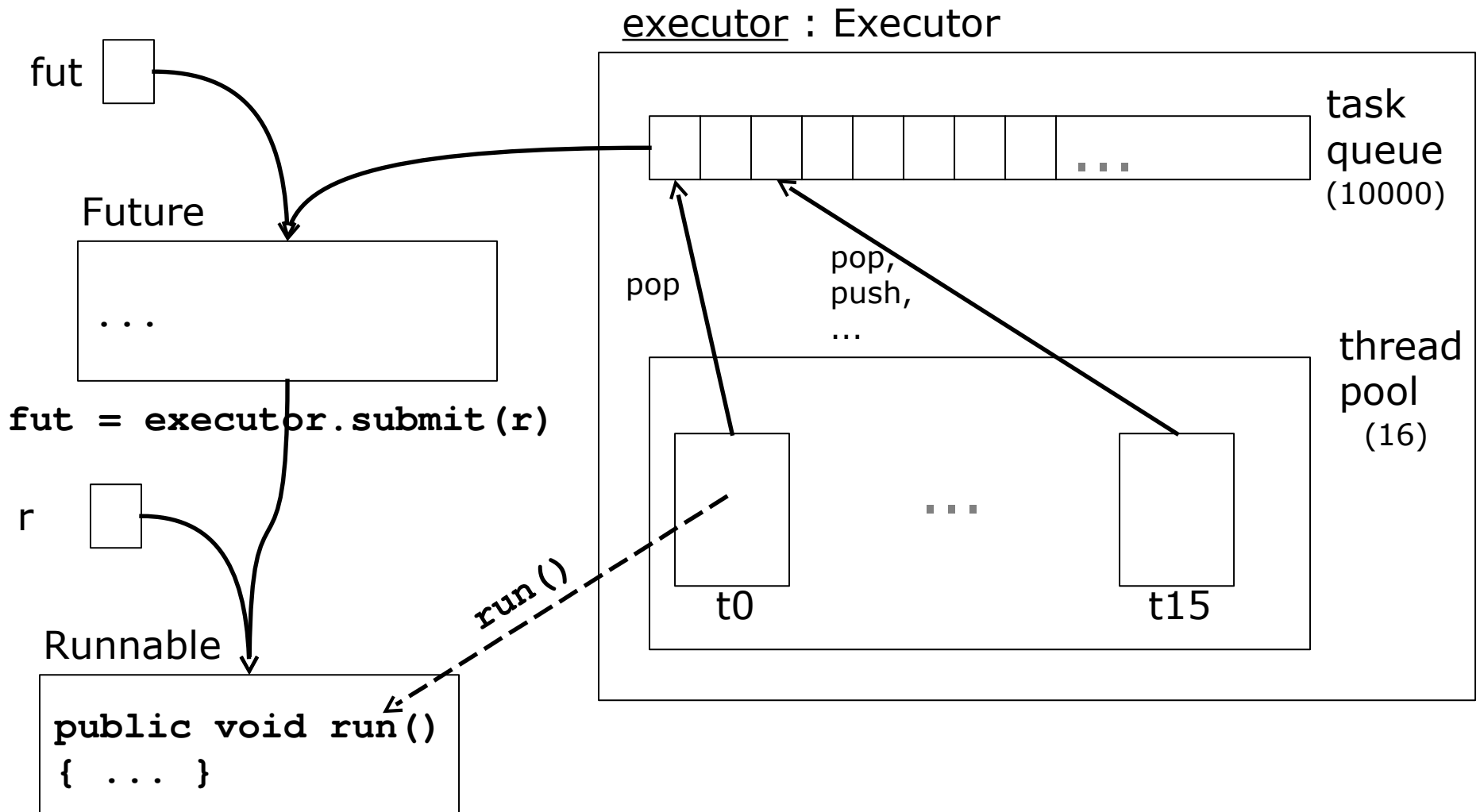
Same, using a lambda

```
Future<?> fut
= executor.submit(() -> { System.out.println("Task ran!"); });
```

- The executor has a pool of threads and uses one of them to run the task
- Use `Future's get ()` to wait for task completion

```
try { fut.get(); }
catch (InterruptedException exn) { System.out.println(exn); }
catch (ExecutionException exn) { throw new RuntimeException(exn); }
```

Dynamics of the executor framework



A task that produces a result

- Make the task from a `Callable<T>`

Future's result type

... same as
Callable's

```
Future<String> fut
= executor.submit(new Callable<String>() {
    public String call() throws IOException {
        return getPage("http://www.wikipedia.org", 10);
    }
});
```

TestTaskSubmit.java

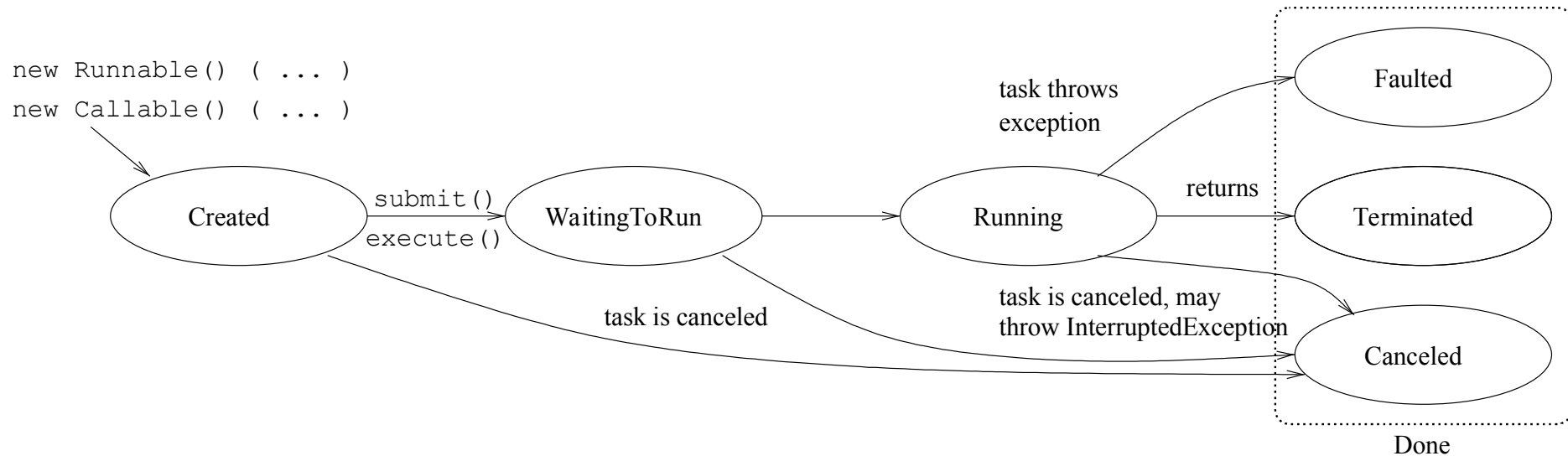
- Use the Future to get the task's result:

```
try {
    String webpage = fut.get();
    System.out.println(webpage);
} catch (InterruptedException exn) { System.out.println(exn); }
catch (ExecutionException exn) { throw new RuntimeExcep...; }
```

Task rules

- Different tasks may run on different threads
 - So objects accessed by tasks must be thread-safe
- A thread running a task can be interrupted
 - So a task can be interrupted
 - So `fut.get()` can throw `InterruptedException`
- Creating a task is fast, takes little memory
- Creating a thread is slow, takes much mem.

The states of a task



- **After `submit` or `execute`**

- a task may be running immediately or much later
- depending on the executor and available threads

Thread creation vs task creation


- Task creation is faster than thread creation

	Thread	Task
Work	6,581 ns	6,612 ns
Create	1,030 ns	77 ns
Create+start/(submit+cancel)	48,929 ns	835 ns
Create+(start/submit)+complete	72,759 ns	21,226 ns

Intel i7 2.4 GHz JVM 1.8

- A task also uses much less memory

Various Java executors

- In class `java.util.concurrent.Executors`:
- `newFixedThreadPool(n)`
 - Fixed number `n` of threads; automatic restart
- `newCachedThreadPool()`
 - Dynamically adapted number of threads, no bound
- `newSingleThreadExecutor()`
 - A single thread; so tasks need not be thread-safe
- `newScheduledThreadPool()`
 - Delayed and periodic tasks; eg clean-up, reporting
- `newWorkStealingPool()`  New in Java 8. Use it
 - Adapts thread pool to number of processors, uses multiple queues; therefore better scalability

Plan for today

- Tasks and the Java executor framework
 - Executors, Runnables, Callables, Futures
- The states of a task
- Task creation overhead
- **Using tasks to count prime numbers**
- Java versus the .NET Task Parallel Library
- Producer-consumer pipelines
- Bounded queues, thread wait and notify
- The states of a thread

Week 1 flashback: counting primes in multiple threads

```
final LongCounter lc = new LongCounter();
Thread[] threads = new Thread[threadCount];
for (int t=0; t<threadCount; t++) {
    final int from = perThread * t,
            to = (t+1==threadCount) ? range : perThread * (t+1);
    threads[t] = new Thread(() -> {
        for (int i=from; i<to; i++)
            if (isPrime(i))
                lc.increment();
    });
}
for (int t=0; t<threadCount; t++)
    threads[t].start();
```

Last thread
has to==range

Thread processes
segment
[from,to)

- Creates one thread for each segment

Counting primes in multiple tasks

```
List<Future<?>> futures = new ArrayList<Future<?>>();
for (int t=0; t<taskCount; t++) {
    final int from = perTask * t,
            to = (t+1 == taskCount) ? range : perTask * (t+1);
    futures.add(executor.submit(() -> {
        for (int i=from; i<to; i++)
            if (isPrime(i))
                lc.increment();
    }));
}
try {
    for (Future<?> fut : futures)
        fut.get();
} catch (...) { ... }
```

Add to shared

Create task, submit to executor, save a future

Wait for all tasks to complete

TestCountPrimesTasks.java

- Creates a task for each segment
- The tasks execute on a thread pool

Tasks with task-local counts

```
List<Callable<Long>> tasks=new ArrayList<Callable<Long>>()
for (int t=0; t<taskCount; t++) {
    final int from = perTask * t,
            to = (t+1 == taskCount) ? range : perTask * (t+1);
    tasks.add(() -> {
        long count = 0;
        for (int i=from; i<to; i++)
            if (isPrime(i))
                count++;
        return count;
    });
}
long result = 0;
try {
    List<Future<Long>> futures = executor.invokeAll(tasks);
    for (Future<Long> fut : futures)
        result += fut.get();
} catch (...) { ... }
```

Create task

Add to
local

Submit tasks, wait for
all to complete, get
futures

Add local task
results

TestCountPrimesTasks.java

Callable<Void> is like Runnable

```
List<Callable<Void>> tasks=new ArrayList<Callable<Void>>()
for (int t=0; t<taskCount; t++) {
    final int from = perTask * t,
        to = (t+1 == taskCount) ? range : perTask * (t+1);
    tasks.add(() -> {
        for (int i=from; i<to; i++)
            if (isPrime(i))
                lc.increment();
        return null;
    });
}
try {
    executor.invokeAll(tasks);
} catch (...) { ... }
```

Create task

Add to shared

Submit tasks, wait for all to complete

TestCountPrimesTasks.java

Type parameters <Void> and <?>

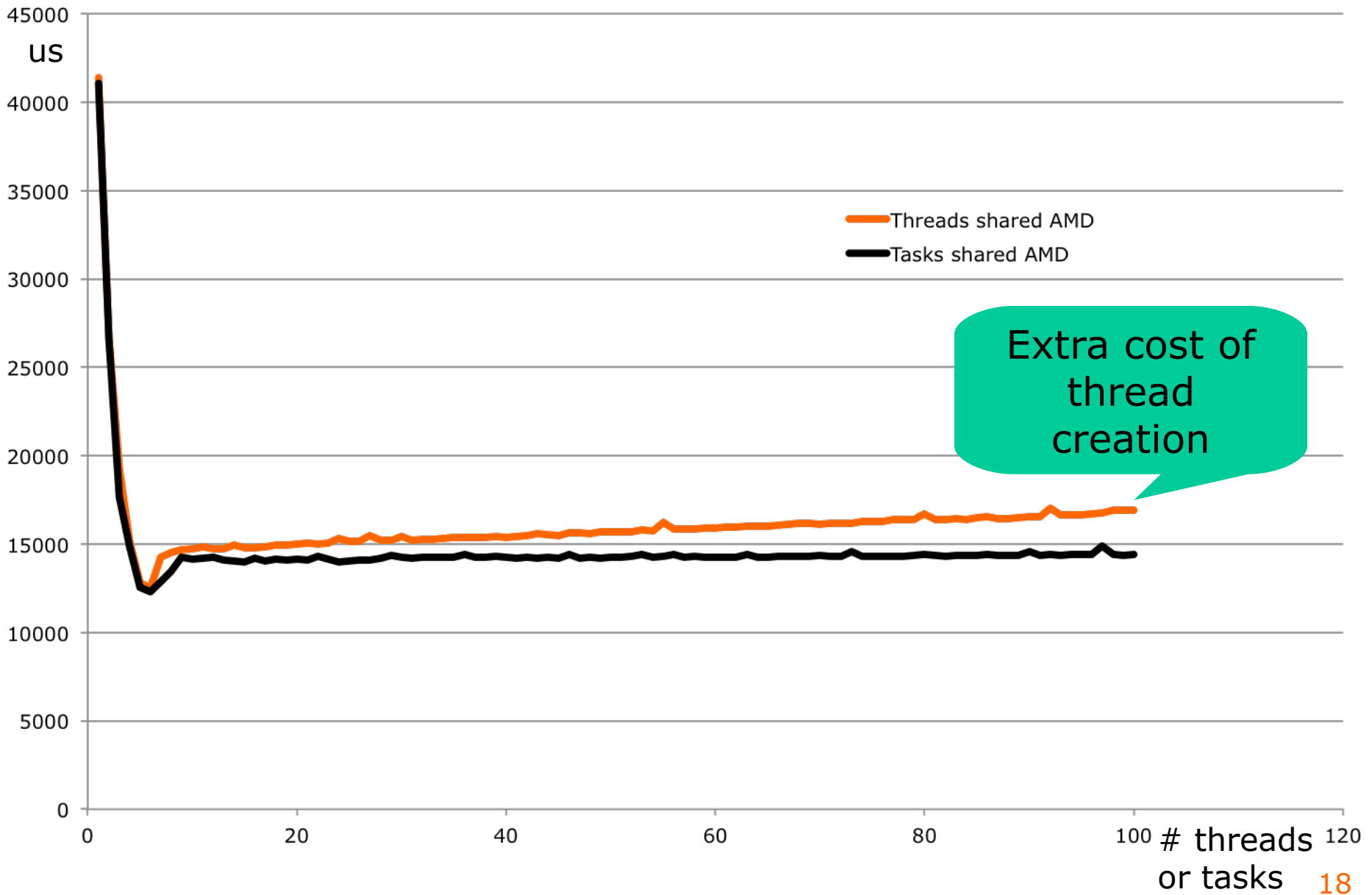
- The type `java.lang.Void` contains only `null`
- `Callable<Void>` requires `Void call() { ... }`
 - Similar to `Runnable`'s `void run() { ... }`
 - With `Future<Void>` the `get ()` returns `null`
- `Future<?>` has an unknown type of value
 - With `Future<?>` the `get ()` returns `null` also
- Java's type system is somewhat muddled
 - Forbids this assignment, so need `Callable<Void>`:

Not
same

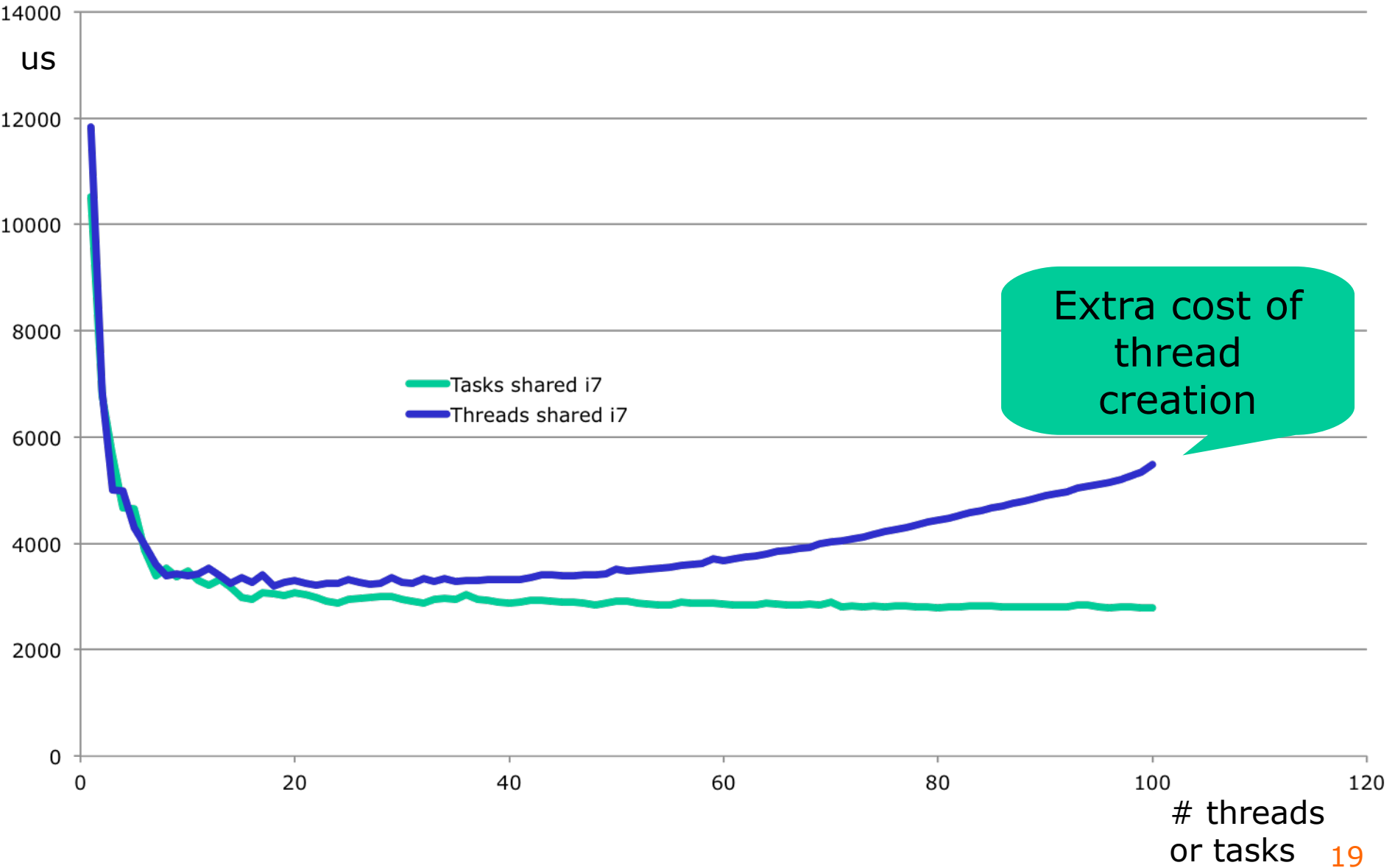
Type `Future<?>`

```
Future<Void> future;  
future = executor.submit(new Runnable() { ... });
```

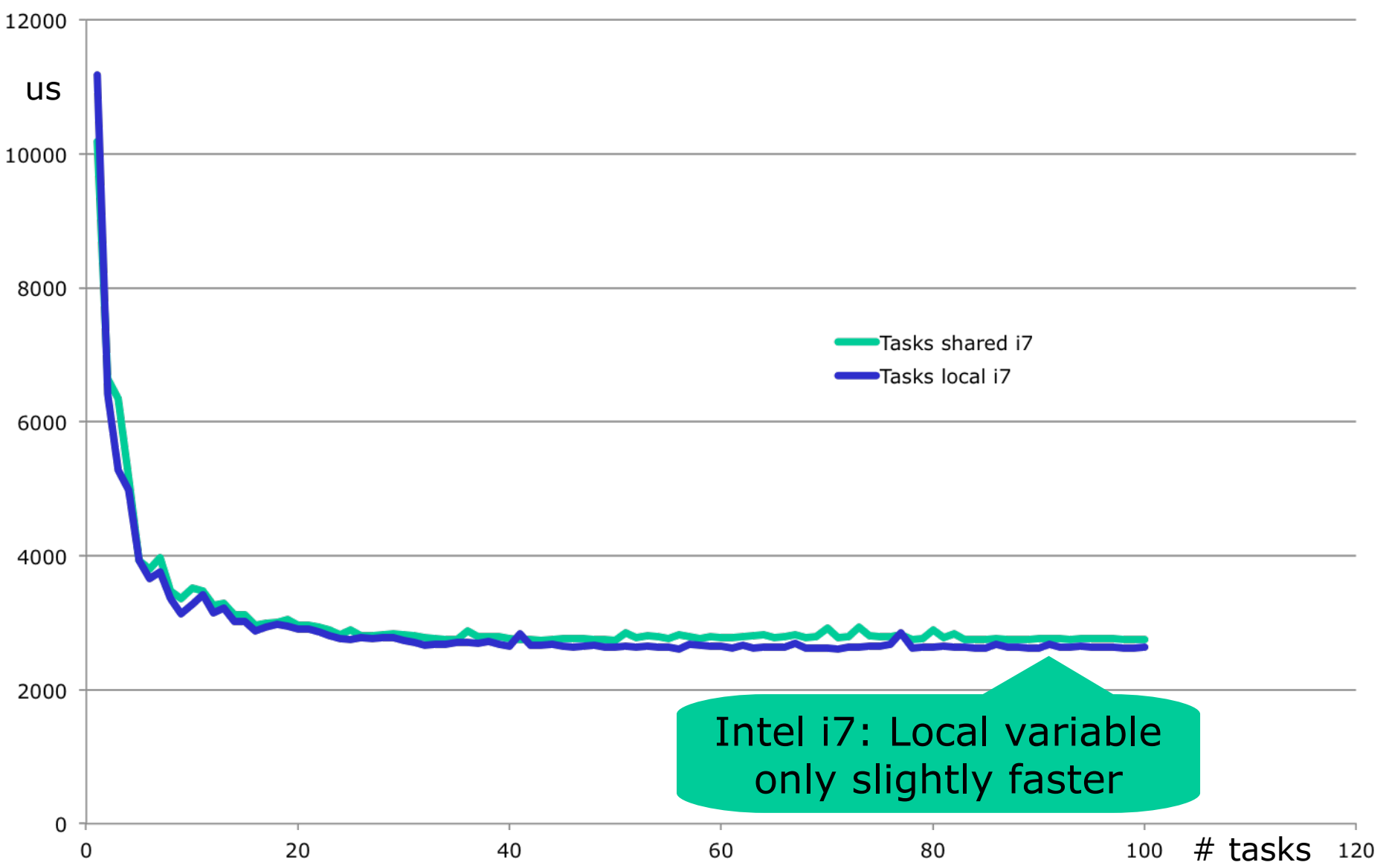
Overhead of creating threads (AMD)



Overhead of creating threads (i7)

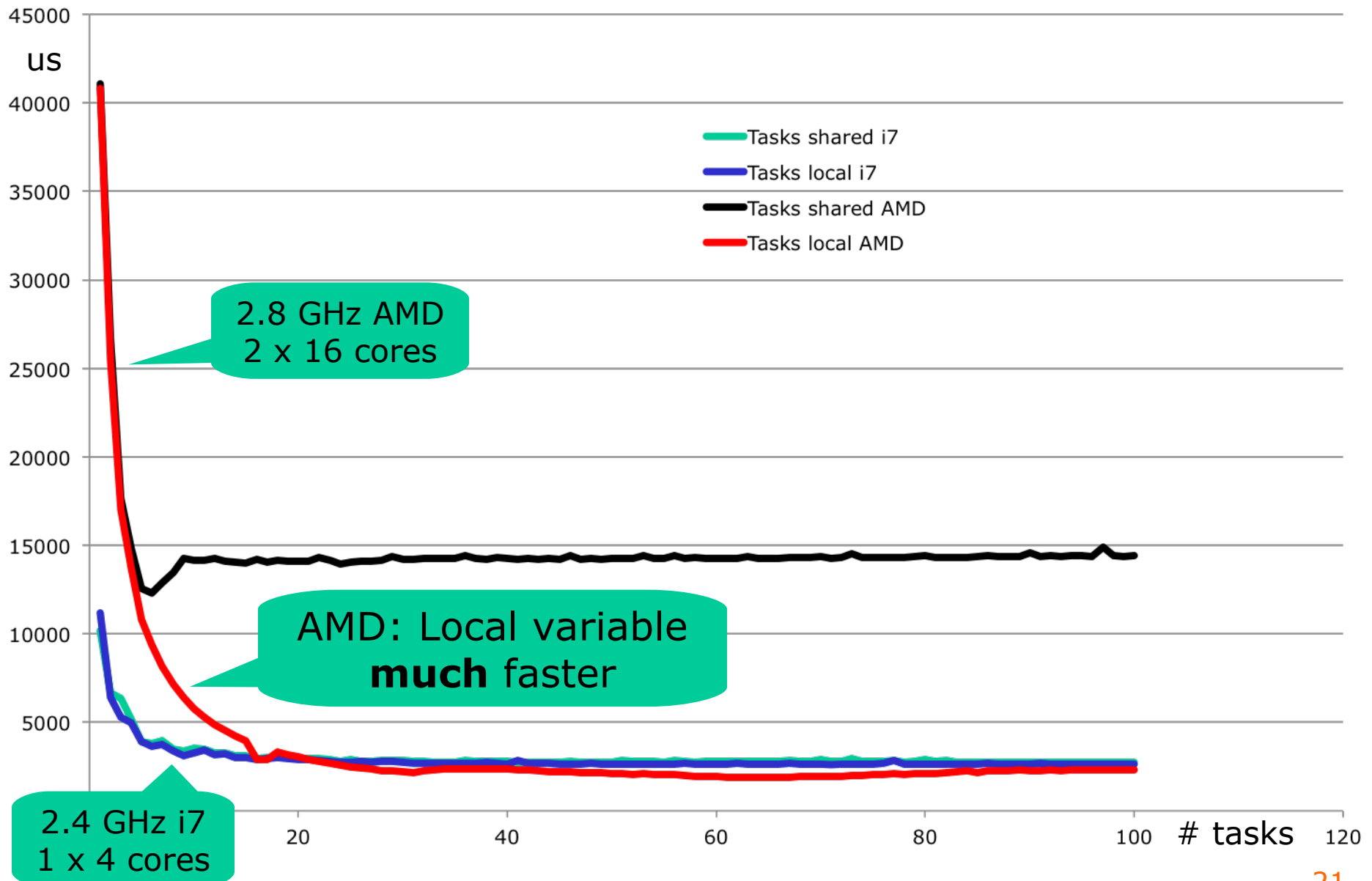


Shared counter vs local counter



Intel i7: Local variable only slightly faster

Computers differ a lot



Plan for today

- Tasks and the Java executor framework
 - Executors, Runnables, Callables, Futures
- The states of a task
- Task creation overhead
- Using tasks to count prime numbers
- **Java versus the .NET Task Parallel Library**
- Producer-consumer pipelines
- Bounded queues, thread wait and notify
- The states of a thread

The .NET Task Parallel Library

- Since C#/.NET 4.0, 2010
- Easier to use and better language integration
 - `async` and `await` keywords in C#
 - .NET class library has more non-blocking methods
 - Java may get non-blocking methods ... sometime
- Namespace `System.Threading.Tasks`
- Class `Task` combines `Runnable` & `Future<?>`
- Class `Task<T>` combines `Callable<T>` and `Future<T>`

- See *C# Precisely* chapters 22 and 23

Parallel prime counts in C#, shared

```
int perTask = range / taskCount;
LongCounter lc = new LongCounter();
Parallel.For(0, taskCount, t =>
{
    int from = perTask * t,
        to = (t+1 == taskCount) ? range : perTask * (t+1);
    for (int i=from; i<to; i++)
        if (isPrime(i))
            lc.increment();
});
return lc.get();
```

Create tasks, submit to standard executor, run

Create task t

- Same concepts as in Java
 - much leaner notation
 - easier to use out of the box
- The tasks are executed on a thread pool
 - in an unknown order

Parallel prime counts in C#, local

```
long[] results = new long[taskCount];
Parallel.For(0, taskCount, t =>
    { int from = perTask * t,
      to = (t+1 == taskCount) ? range : perTask * (t+1);
      long count = 0;
      for (int i=from; i<to; i++)
          if (isPrime(i))
              count++;
      results[t] = count;
    });
return results.Sum();
```

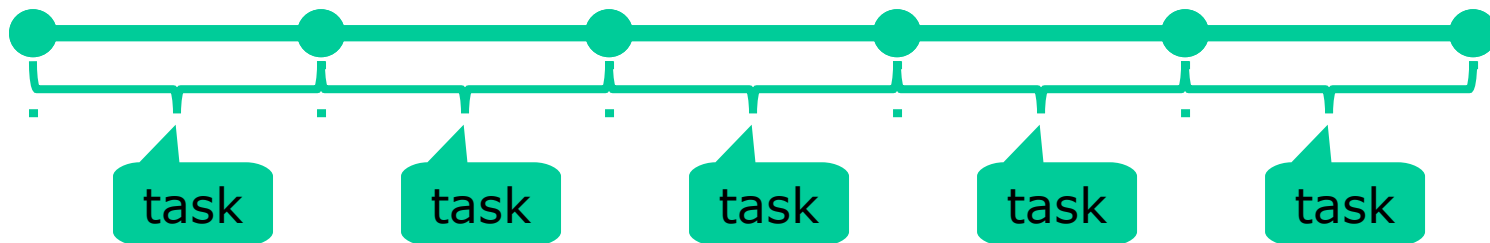
- Q: Why safe to write to **results** array?

Plan for today

- Tasks and the Java executor framework
 - Executors, Runnable, Callable, Future
- The states of a task
- Task creation overhead
- Using tasks to count prime numbers
- Java versus the .NET Task Parallel Library
- **Producer-consumer pipelines**
- Bounded queues, thread wait and notify
- The states of a thread

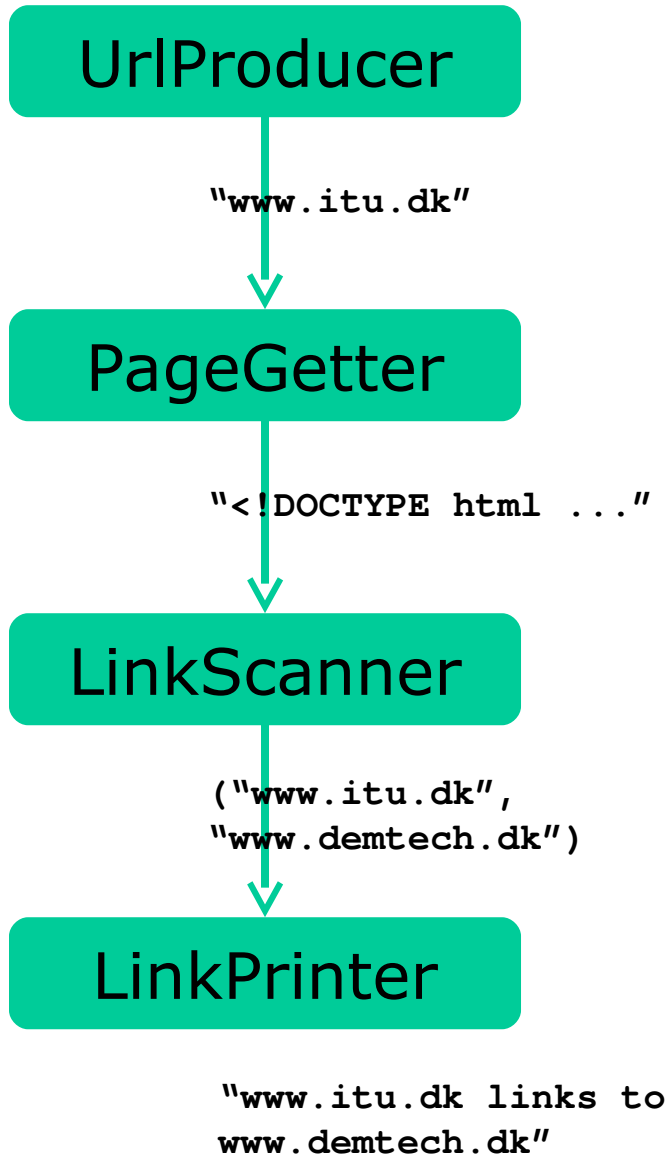
Concurrent pipelines (Goetz §5.3)

- We parallelized prime counting by splitting the work into chunks:



- A different way is to create a *pipeline*
- Example problem: Given long list of URLs,
 - For each URL,
 - download the webpage at that URL
 - scan the webpage for links ` ...`
 - for each link, print `"url links to link"`

Pipeline to produce URL, get webpage, scan for links, and print them



- There are four stages
- They can run in parallel
 - On four threads
 - Or as four tasks
- Each does a simple job
- Two stages communicate via a blocking queue
 - `queue.put(item)` sends data item to next stage; blocks until room for data
 - `queue.take()` gets data item from previous stage; blocks until data available

Sketch of a one-item queue

```
interface BlockingQueue<T> {  
    void put(T item);  
    T take();  
}
```

```
class OneItemQueue<T> implements BlockingQueue<T> {  
    private T item;  
    private boolean full = false;  
    public void put(T item) {  
        synchronized (this) {  
            full = true;  
            this.item = item;  
        }  
    }  
    public T take() {  
        synchronized (this) {  
            full = false;  
            return item;  
        }  
    }  
}
```

Java monitor pattern, good

But: what if already full?

If queue full, we must wait for **another** thread to **take()** first

But: What if queue empty?

Other thread can **take()** only if we release lock first

Useless

Using wait() and notifyAll()

```
public void put(T item) {
    synchronized (this) {
        while (full) {
            try { this.wait(); }
            catch (InterruptedException exn) { }
        }
        full = true;
        this.item = item;
        this.notifyAll();
    }
}
```

If queue full, wait for notify from other thread

When non-full, save item, notify all waiting threads

- **this.wait()**: release lock on **this**; do nothing until notified, then acquire lock and continue
 - Must hold lock on **this** before call
- **this.notifyAll()**: tell all threads **wait()**ing on **this** to wake up
 - Must hold lock on **this**, and keeps holding it

The take() method is similar

```
public T take() {  
    synchronized (this) {  
        while (!full) {  
            try { this.wait(); }  
            catch (InterruptedException exn) { }  
        }  
        full = false;  
        this.notifyAll();  
        return item;  
    }  
}
```

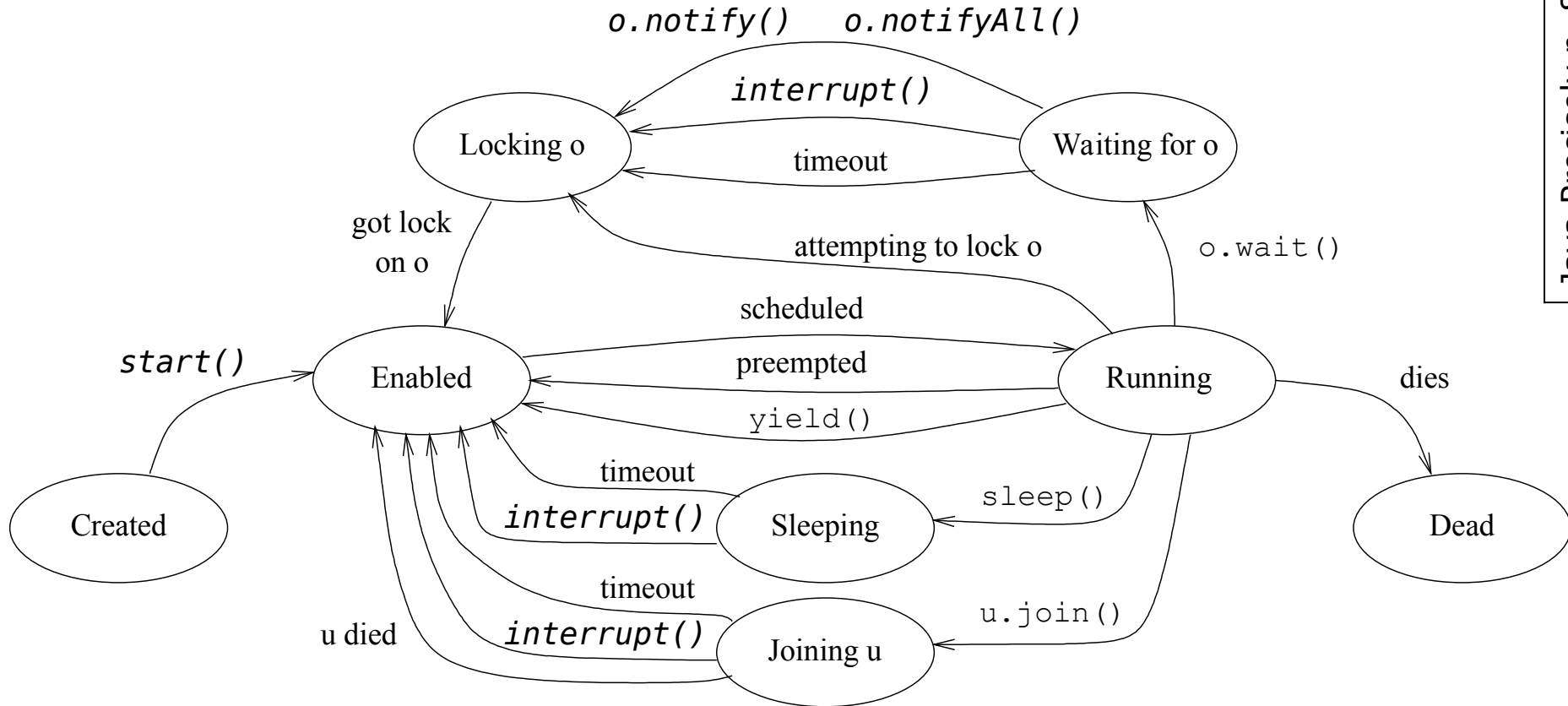
If queue empty, wait for notify from other thread

When non-empty, take item, notify all waiting threads

- Only works if **all** methods locking on the queue are written correctly
- MUST do the wait() in a while loop; Q: Why?

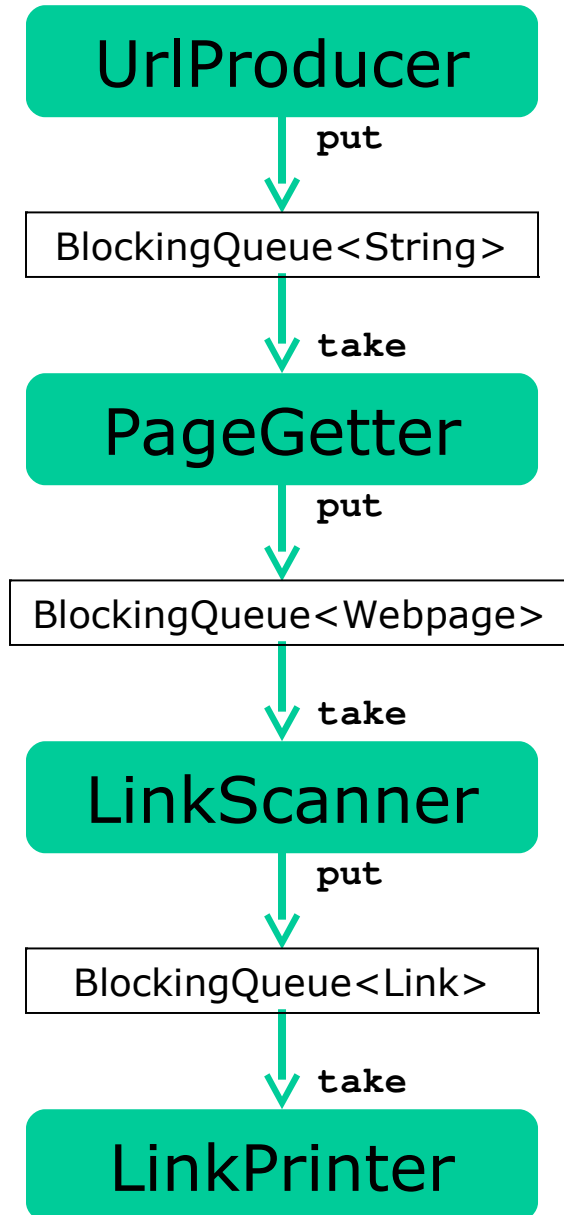
Always use the wait loop idiom to invoke the wait method; never invoke it outside of a loop. The loop serves to test the condition before and after waiting.

Java Thread states



- `o.wait()` is an action of the running thread itself
- `o.notify()` is an action by another thread, on the waiting one
- `scheduled`, `preempted`, ... are actions of the system

Producer-consumer pattern: Pipeline stages and connecting queues



- The first stage is a producer only
- The middle stages are both consumers and producers
- The last stage is only a consumer

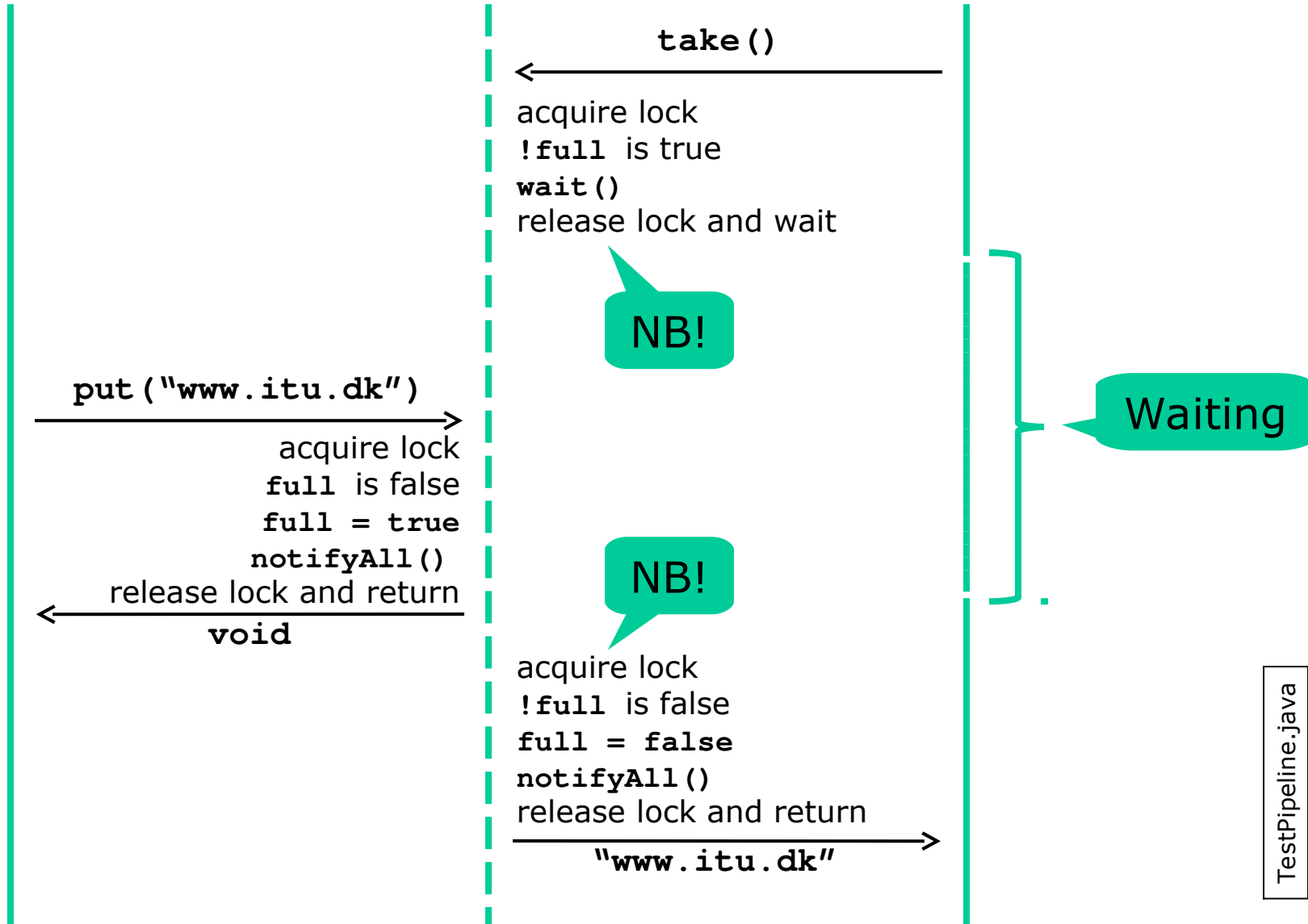
- A queue connects producer(s) to consumer(s) in a thread-safe way

How wait and notifyAll collaborate

UrlProducer
(active thread)

OneItemQueue
(passive object)

PageGetter
(active thread)



TestPipeline.java

Stages 1 and 2

```
class UrlProducer implements Runnable {
    private final BlockingQueue<String> output;
    public UrlProducer(BlockingQueue<String> output) {
        this.output = output;
    }
    public void run() {
        for (int i=0; i<urls.length; i++)
            output.put(urls[i]);
    }
}
```

Produce URLs

```
class PageGetter implements Runnable {
    ...
    public void run() {
        while (true) {
            String url = input.take();
            try {
                String contents = getPage(url, 200);
                output.put(new Webpage(url, contents));
            } catch (IOException exn) { System.out.println(exn); }
        }
    }
}
```

Transform
URL to
webpage

Stages 3 and 4

```
class LinkScanner implements Runnable {
    ...
    private final static Pattern urlPattern
        = Pattern.compile("a href=\"(\\p{Graph}*)\"")
    public void run() {
        while (true) {
            Webpage page = input.take();
            Matcher urlMatcher = urlPattern.matcher(page.contents);
            while (urlMatcher.find()) {
                String link = urlMatcher.group(1);
                output.put(new Link(page.url, link));
            }
        }
    }
}
```

Transform
web page to
link "stream"

```
class LinkPrinter implements Runnable {
    ...
    public void run() {
        while (true) {
            Link p = input.take();
            System.out.printf("%s links to %s%n", p.from, p.to);
        }
    }
}
```

Consume links
and print
them

Putting stages and queues together

```
final BlockingQueue<String> urls = new OneItemQueue<String>();  
final BlockingQueue<Webpage> pages = new OneItemQueue<Webpage>();  
final BlockingQueue<Link> refPairs = new OneItemQueue<Link>();  
Thread t1 = new Thread(new UrlProducer(urls));  
Thread t2 = new Thread(new PageGetter(urls, pages));  
Thread t3 = new Thread(new LinkScanner(pages, refPairs));  
Thread t4 = new Thread(new LinkPrinter(refPairs));  
t1.start(); t2.start(); t3.start(); t4.start();
```

TestPipeline.java

- Each stage does *one* job
 - Simple to implement and easy to modify
 - Separation of concerns, simple control flow
- Easy to add new stages
 - For instance, discard duplicate links
- Can achieve high throughput
 - May run multiple copies of a slow stage

“Prefer concurrency utilities to wait and notify”

Bloch item 69

- It's instructive to use `wait` and `notify`
- ... but easy to do it wrong
- Package `java.util.concurrent` has
 - `BlockingQueue<T>` interface
 - `ArrayBlockingQueue<T>` class and much more
- Better use those in practice

- ... or make a pipeline with Java 8 streams
 - Simpler, and *very* easy to parallelize

Using Java 8 streams instead

- Package `java.util.stream`
- A `Stream<T>` is a source of T values
 - Lazily generated
 - Can be transformed with `map(f)` and `flatMap(f)`
 - Can be filtered with `filter(p)`
 - Can be consumed by `forEach(action)`
- Generally simpler than concurrent pipeline

```
Stream<String> urlStream
    = Stream.of(urls);
Stream<Webpage> pageStream
    = urlStream.flatMap(url -> makeWebPageOrNone(url, 200));
Stream<Link> linkStream
    = pageStream.flatMap(page -> makeLinks(page));
linkStream.forEach(link ->
    System.out.printf("%s links to %s%n", link.from, link.to));
```

Making the stages run in parallel

```
Stream<String> urlStream
    = Stream.of(urls).parallel();
Stream<Webpage> pageStream
    = urlStream.flatMap(url -> makeWebPageOrNone(url, 200));
Stream<Link> linkStream
    = pageStream.flatMap(page -> makeLinks(page));
linkStream.forEach(link ->
    System.out.printf("%s links to %s\n", link.from, link.to));
```

TestStreams.java

- Magic? No!
- Divides streams into substream chunks
- Evaluates the chunks in tasks
- Runs tasks on an executor called ForkJoinPool
 - Using a thread pool and work stealing queues
 - More precisely ForkJoinPool.commonPool()

So easy. Why learn about threads?

- Parallel streams use tasks, run on threads
- Should be **side effect free** and take no locks
- Otherwise all the usual thread problems:
 - updates must be made atomic (by locking)
 - updates must be made visible (by locking, volatile)
 - deadlock risk if locks are taken

Side-effects

Side-effects in behavioral parameters to stream operations are, in general, **discouraged**, as they can often lead to unwitting violations of the statelessness requirement, as well as other thread-safety hazards.

If the behavioral parameters do have side-effects, unless explicitly stated, there are **no guarantees as to the visibility of those side-effects** to other threads, nor are there any guarantees that different operations on the "same" element within the same stream pipeline are executed in the same thread. Further, the ordering of those effects may be surprising.

How are Java streams implemented?

- Splititerator = splittable iterator

```
interface Splititerator<T> {  
    void forEachRemaining(Consumer<T> action);  
    boolean tryAdvance(Consumer<T> action);  
    void Splititerator<T> trySplit();  
}
```

- Many method calls (well inlined/fused by the JIT)
- Parallelization
 - Divide stream into chunks
 - Process each chunk in a task (Haskell “spark”)
 - Run on thread pool using work-stealing queues
 - ... thus similar to Haskell parBuffer/parListChunk

Example: array spliterator def.

```
class ArraySpliterator<T> implements Spliterator<T> {  
    private final Object[] array; // underlying array  
    private int index;           // next index, modified on advance/split  
    private final int fence;     // one past last index
```

```
    public boolean tryAdvance(Consumer<? super T> action) {  
        if (index >= 0 && index < fence) {  
            action.accept((T) array[index++]);  
            return true;  
        } else  
            return false;  
    }
```

Consume
one element

```
    public Spliterator<T> trySplit() {  
        int lo = index, mid = (lo + fence) >>> 1;  
        return (lo >= mid)  
            ? null  
            : new ArraySpliterator<>(array, lo, index = mid, characteristics);  
    }
```

Split into two
Spliterators

```
    public void forEachRemaining(Consumer<? super T> action) { ... }  
}
```

Example: array splitter

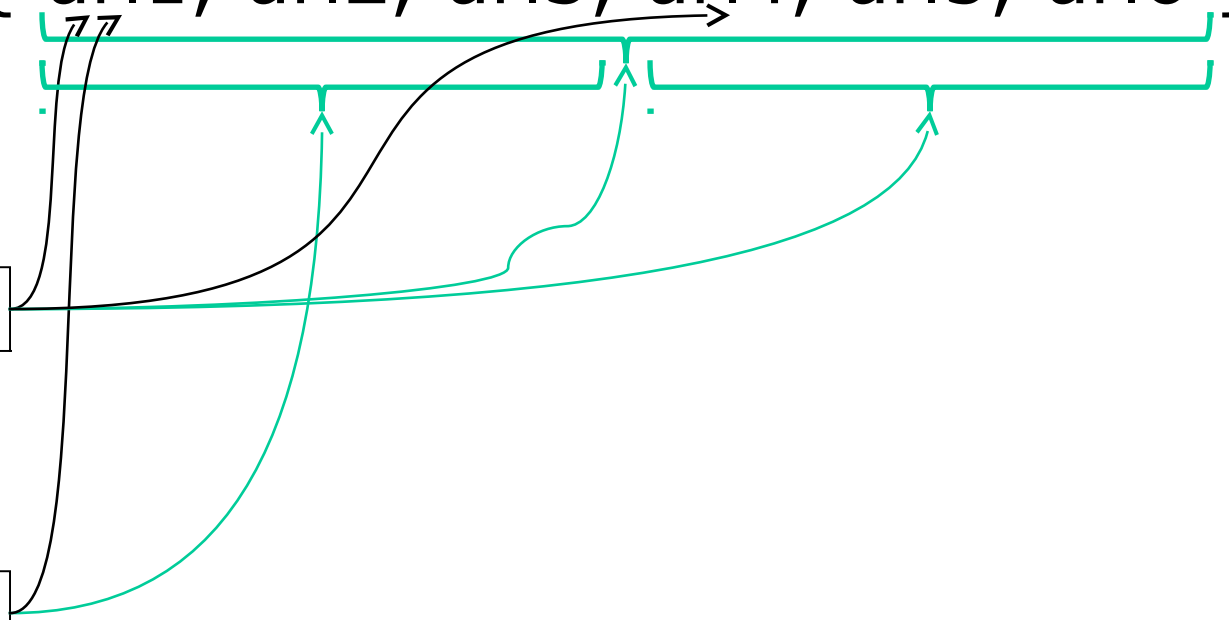
- array arr = { url1, url2, url3, url4, url5, url6 }

s1 = Stream.of(arr)
.splitter()

ArraySplitter s1

s2 = s1.trySplit()

ArraySplitter s2



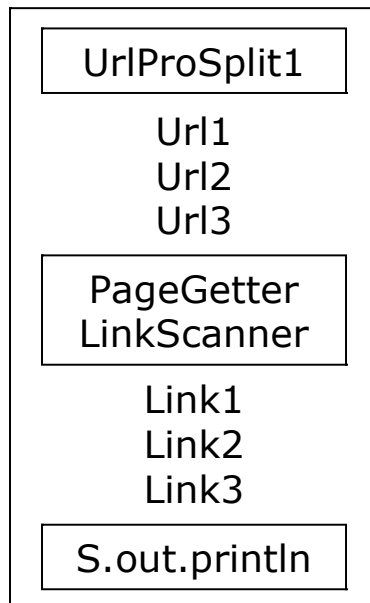
Thread vs Streams parallel pipelines

UrlProducer → PageGetter → LinkScanner → LinkPrinter

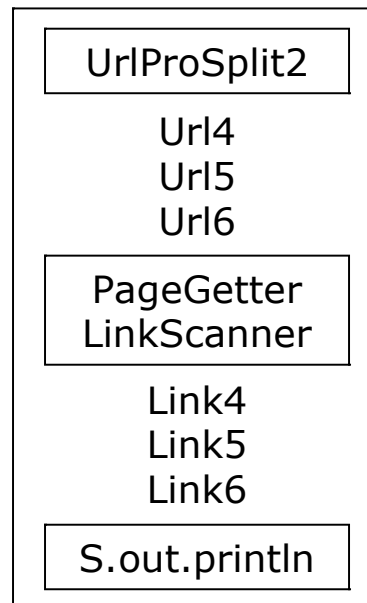
Time ↓

Url1		
Url2	Page1	
Url3	Page2	Link1
Url4	Page3	Link2
Url5	Page4	Link3
Url6	Page5	Link4

ForEachTask1



ForEachTask2



This week

- Reading
 - Goetz et al chapters 5.3, 6 and 8
 - Bloch items 68, 69
- Hand-in week 5
 - Show that you can use tasks and the executor framework, and modify a concurrent pipeline
- Read before next week's lecture
 - Goetz chapter 11, 13.5