



Real-Time Java[®] for Ada Programmers

Ben Brosgol
brosgol@adacore.com

Tutorial SF3 – SIGAda 2004
Atlanta, GA
USA

14 November 2004

US Headquarters:

104 Fifth Avenue, 15th Floor
New York, NY 10011

+1-212-620-7300 (voice)
+1-212-807-0162 (FAX)

www.adacore.com

Part 1 - Basics

- Language requirements for real-time programming
- Java overview, thread model, and assessment for real-time applications
- Summary of underlying issues: priority inversion, garbage collection

Part 2 – Real-Time Specification for Java (“RTSJ”)

- Overview and design goals
- Scheduling
- Memory Management

Part 3 – RTSJ (continued)

- Synchronization / priority inversion management
- Asynchronous event handling

Part 4 – RTSJ (concluded) and other real-time Java topics

- Asynchronous transfer of control
- Time and timers
- Low-level features
- Real-time Java in practice
- Conclusions

Audience background

- Some knowledge of (sequential) Java
- Some knowledge of real-time issues
- No knowledge of Java thread model
- No knowledge of the real-time Java proposals

Presenter background

- Primary member of “Expert Group” for original design of RTSJ
- Member of Technical Interpretation Committee for RTSJ
- Reviewer of J-Consortium Core Extensions
- Author of several Ada / Java comparison papers
- Frequent presenter at Ada Europe and SIGAda conferences
- Participant in Ada 83 and Ada 95 design efforts
- Member of AdaCore senior technical staff

General requirements

- Language features promoting safety/reliability
 - Compile-time checking
 - Encapsulation
 - Deterministic language semantics
- Predictability
 - In deadline-sensitive parts of the program, need to know worst-case execution time for language constructs, API
- Performance (?)

Concurrency

- Adequate range of priority values
- Well-defined semantics for scheduling
- Safe / efficient mutual exclusion, including “state notification”
- Safe / efficient coordination / communication
- Avoidance of unbounded priority inversion
- Support for common real-time idioms such as:
 - Periodic activities
 - Various kinds of events / semaphores

Memory management

- Cannot run out of storage or fragment
- Safe/predictable storage reclamation

Asynchrony

- Asynchronous events / event handlers
 - Connection with interrupts
- Asynchronous Transfer of Control
 - Timeout
 - Thread termination

Time

- Support for high-resolution time (millis and nanos), both absolute and relative
- Support for various kinds of timers, clocks

Low-level support

- Access to hardware-specific features

API

- Functionality, "thread safety"

Java language

- “Pure” Object-Oriented language in the style of Smalltalk
 - Single inheritance of classes, “multiple inheritance” of “interfaces”
- Built-in support for exception handling, threads
- Well-defined semantics, at least for sequential features
- Flexible, dynamic language
 - Classes are run-time objects
 - All non-primitive data go on the heap
- Large emphasis on safety, security (downloadable “applets”)
- Garbage collection required

Java Virtual Machine

- Portable, interpretable binary format for Java classes (class file)
- Instruction set tailored to Java language

Java API

- Extensive set of “packages” for a wide variety of application domains
- “Core” libraries, networking, I/O, collections, database, reflection, ...

Unit of concurrency is the *thread*

- An instance of the class `java.lang.Thread` or one of its subclasses
- `run()` method = algorithm performed by each instance

Either extend `Thread`, or implement the `Runnable` interface

- Override/implement `run()`
- All threads are dynamically allocated

Example

```
public class Writer extends Thread{
    final int count;
    public Writer(int count){this.count=count;}

    public void run(){
        for (int i=1; i<=count; i++){
            System.out.println("Hello " + i);
        }
    }
    public static void main( String[] args ) {
        Writer w = new Writer(60);
        w.start(); // New thread of control invokes w.run()
        System.out.println( "In main thread" ); // May be concurrent with w
    } // May return while w is still running
}
```

Lifetime properties

- Constructing a thread creates thread resources (stack, etc.)
- “Activation” of thread `t` is explicit, by invoking `t.start()`
- Started thread runs “concurrently” with parent
- Thread terminates when its `run()` method returns
- Parent does not need to wait for children to terminate

Mutual exclusion

- Shared data (volatile fields)
- synchronized blocks/methods

Thread coordination/communication

- Pass data to new thread via constructor
- Pulsed event – `obj.wait()` / `obj.notify()`
- Broadcast event – `obj.wait()` / `obj.notifyAll()`
- `t.join()` suspends caller until the target thread `t` completes



Time

- `sleep(millis)`
- Timeout on `obj.wait(millis)`

Scheduling/priorities

- Priority is in range 1..10
- Thread can change/interrogate its own or another thread's priority
- `yield()` gives up the processor

Asynchrony

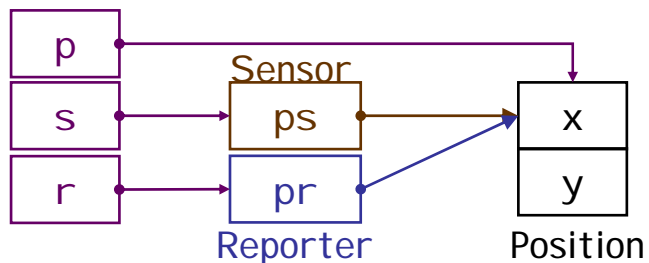
- `interrupt()` sets a bit that can be polled
- Asynchronous termination
 - `stop()` throws asynchronous exception (deprecated, dangerous)
 - `destroy()` kills a thread (unimplemented, dangerous)
- `suspend()` / `resume()` have been deprecated
- No asynchronous exceptions
- Thread propagating an unhandled exception terminates silently, but first calls its `ThreadGroup`'s `uncaughtException()` method

Other functionality

- `ThreadGroup` class
 - Supplies `uncaughtException` method
- Daemon threads
 - Allows application to terminate though such threads are still running

A Position object is periodically updated by a Sensor object and is periodically displayed by a Reporter object

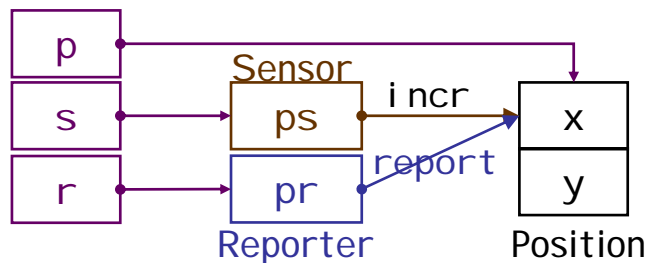
```
class Position{
    double x, y;
}
```



```
class Test{
    public static void main(String[] args){
        Position p = new Position();
        Sensor s = new Sensor( p );
        Reporter r = new Reporter( p );
        s.start();
        r.start();
    }
}
```

```
class Sensor extends Thread{
    private final Position ps;
    Sensor( Position p ){
        ps=p;
    }
    public void run(){
        while(...){
            synchronized( ps ){
                ps.x += ...;
                ps.y += ...;
            }
            sleep(...);
        }
    }
}
```

```
class Reporter extends Thread{
    private final Position pr;
    Reporter( Position p ){
        pr=p;
    }
    public void run(){
        while(...){
            synchronized( pr ){
                System.out.println( pr.x );
                System.out.println( pr.y );
            }
            sleep(...);
        }
    }
}
```



```

class Position{
    private double x, y;
    public synchronized void incr(double xi nc,
                                   double yi nc){
        x += xi nc;
        y += yi nc;
    }
    public synchronized void report() {
        System.out.println(x);
        System.out.println(y);
    }
}
  
```

```

class Sensor extends Thread{
    private final Position ps;
    Sensor( Position p ){
        ps=p;
    }
    public void run(){
        double xi nc, yi nc;
        while(...){
            xi nc = ...;
            yi nc = ...
            ps.incr(xi nc, yi nc);
            sleep(...);
        }
    }
}
  
```

```

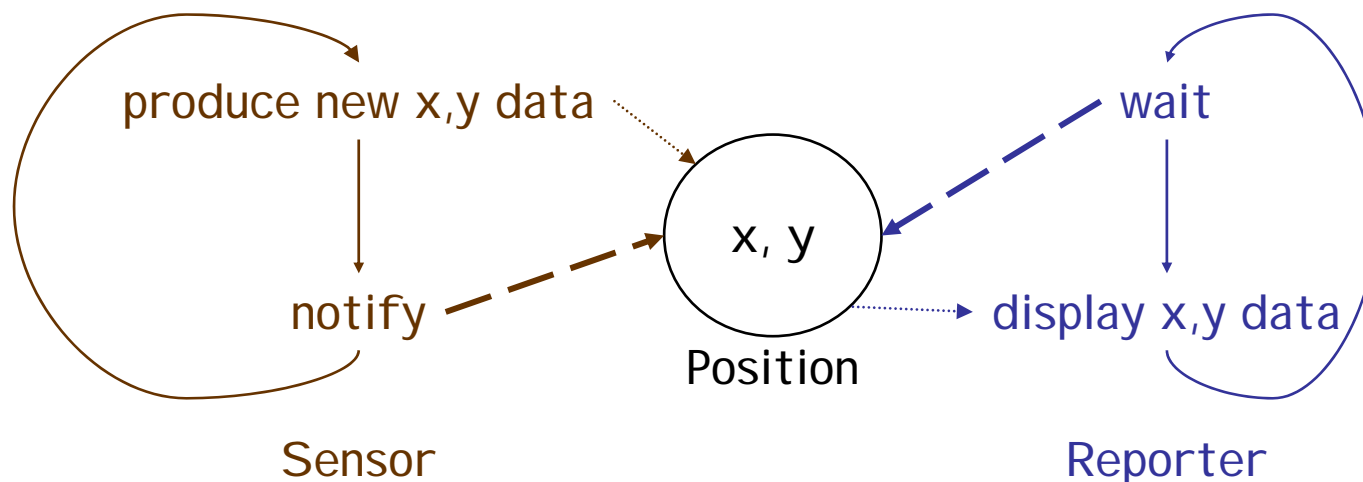
class Reporter extends Thread{
    private final Position pr;
    Reporter( Position p ){
        pr=p;
    }
    public void run(){
        while(...){
            pr.report();
            sleep(...);
        }
    }
}
  
```

A Reporter object may display the same x, y pair several times

Variation: “no stale data”

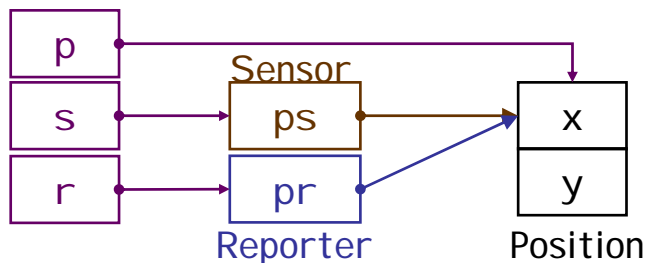
- Ensure that a Reporter always obtains an x, y pair that is set *after* the Reporter requests the data
- Do not care if some x, y pairs go unread
 - Sensor object may generate new values more frequently than the Reporter is retrieving them
 - No need to have a Sensor object wait for the new x, y values to be retrieved, or to buffer them

Schematic



Coordination via wait() and notify() (2)

```
class Position{
  double x, y;
}
```



```
class Test{
  public static void main(String[] args){
    Position p = new Position();
    Sensor s = new Sensor( p );
    Reporter r = new Reporter( p );
    s.start();
    r.start();
  }
}
```

```
class Sensor extends Thread{
  private final Position ps;
  Sensor( Position p ){
    ps=p;
  }
  public void run(){
    while(...){
      synchronized( ps ){
        ps.x += ...;
        ps.y += ...;
        ps.notify();
      }
      sleep(...);
    }
  }
}
```

```
class Reporter extends Thread{
  private final Position pr;
  Reporter( Position p ){
    pr=p;
  }
  public void run(){
    while(...){
      synchronized( ps ){
        pr.wait();
        System.out.println( pr.x );
        System.out.println( pr.y );
      }
      sleep(...);
    }
  }
}
```

Goals

- Provide high-level, high-performance concurrency building blocks
- Replace many uses of wait / notify

New APIs, based on Doug Lea's util.concurrent library

- Mutual exclusion support
 - Explicit locks
 - Explicit condition variables
 - Atomic variables for thread-safe operations with no synchronization
- Communication idioms
 - Multi-way synchronization (rendezvous)
 - Semaphores
- Thread-aware collection classes
- Thread pools
- Thread schedulers
- Nanosecond timing

Part of J2SE 5.0 (Java 1.5)

- Developed independently of Real-Time Specification for Java
- Contained in java.util.concurrent package

General benefits

- Language security and (in general) well-defined semantics
- Portability at multiple levels (“Write Once, Run Anywhere”)
- Extensive API

Technical features / expressiveness / flexibility

- Support for software engineering (encapsulation, OOP, exceptions...)
- Built-in feature for concurrency (threads)
- Dynamic loading attractive in some segments such as telecom

Advantages over other languages

- Safer than C
- Simpler than C++
- More popular than Ada

Pragmatics / politics

- Organization adopting Java as an “enterprise” language wish to use Java for real-time

Error-prone

- Requires cooperation by the accessing threads
 - Even if all methods are synchronized, an errant thread can access non-private fields without synchronization
- Subtle bug: constructor or synchronized instance method making non-synchronized access to static field
- “Nested monitor” problem

Subtleties in practice

- Not always clear when a method needs to be declared as synchronized
- Complex interactions with other features (e.g. when are locks released)
- Locking is hard to get right (exacerbated by absence of nested objects)

Effect not always clear from source syntax

- A non-synchronized method may be safe to invoke from multiple threads
- A synchronized method might not be safe to invoke from multiple threads

Thread communication/synchronization issues

- `wait()` and `notify()/notifyAll()` are low-level constructs that must be used very carefully
 - `"while (!condition) {obj.wait()}"` needed
- Limited mechanisms for direct inter-thread communication
- Synchronized code that changes object's state must explicitly invoke `notify()` or `notifyAll()`
- No syntactic distinction between signatures of synchronized method that may suspend a caller and one that does not
- Only one wait set per object (versus per associated "condition")

Public thread interface issues

- The need to explicitly initiate a thread by invoking its `start()` method allows several kinds of programming errors
- Although `run()` is part of a thread class's public interface, invoking it explicitly is generally an error

"Inheritance Anomaly"

- Subclass may need to redo superclass synchronization logic

Lack of some features useful for software engineering

- Operator overloading, strongly typed primitive types, ...

Thread model deficiencies

- Priority range (1..10) too narrow
- Priority semantics are implementation dependent and fail to prevent unbounded priority inversion
- Relative sleep() not sufficient for periodicity

Memory management unpredictability

- Predictable, efficient garbage collection appropriate for real-time applications is not (yet) in the mainstream
- Java lacks stack-based objects (arrays and class instances)
- Heap used for exceptions thrown implicitly as an effect of other operations

Asynchrony deficiencies

- Event handling requires dedicated thread
- interrupt() not sufficient
- stop() and destroy() deprecated or dangerous or both

Section 17.12 of the Java Language Specification

- “Every thread has a *priority*. When there is competition for processing resources, threads with higher priority are *generally* executed in preference to threads with lower priority. *Such preference is not, however, a guarantee that the highest priority thread will always be running, and thread priorities cannot be used to reliably implement mutual exclusion.*”

Problems for real-time applications

- Impossible to guarantee that deadlines will be met for periodic threads
 - May get priority inversion
- No guarantee that priority is used for selecting a thread to unblock when a lock is released
- No guarantee that priority is used for selecting which thread is awakened by a `notify()`, or which thread awakened by `notifyAll()` is selected to run

Run-time issues

- Dynamic class loading is expensive, not easy to see when it will occur
- Array initializers ⇒ run-time code

OOP has not been embraced by the real-time community

- Dynamic binding complicates analyzability
- Garbage Collection defeats predictability

Lack of features for accessing the underlying hardware

Performance questions

API Issues

- “Standard” API would need to be rewritten for predictability
- A class’s “interface” is more than its public and protected members
 - In general it includes some implementation characteristics
 - E.g. does it allocate objects, can it block

Some JVM opcodes require non-constant amount of time

```
class Periodic extends Thread{
    static interface Action{ void doIt(); }
    private final Action action;
    private final long periodMillis;
    Periodic(String name, long periodMillis, int priority, Action action){
        super(name);
        this.periodMillis = periodMillis;
        this.action = action;
        this.setPriority(priority);
    }
    public void run(){
        long nextTime = System.currentTimeMillis();
        while(true){
            this.action.doIt(); // action to be performed
            nextTime += periodMillis;
            try{
                sleep(nextTime - System.currentTimeMillis());
            } catch (InterruptedException e){ return; }
        }
    }
    public static void main(String[] args){
        Periodic p1 = new Periodic("High", 100, 10, new Action(){...});
        Periodic p2 = new Periodic("Low", 500, 2, new Action(){...});
        p1.start(); p2.start(); ...
        p1.interrupt(); p2.interrupt();
    }
}
```

Problems:

- No guarantee that priorities will be obeyed
- Relative delay in sleep may lead to missed deadlines

java.util.TimerTask

- A Runnable object whose run() method is invoked, either once or repeatedly, by a Timer

java.util.Timer

- An object with an associated background thread that invokes the Timer's TimerTasks

```
import java.util.*;
class PeriodicExample{
    public static void main(String[] args){
        Timer periodic1 = new Timer();
        Timer periodic2 = new Timer();
        periodic1.scheduleAtFixedRate(
            new TimerTask(){
                public void run(){...} // periodic action 1
            }, 0, 100); // start now, period = 100ms
        periodic2.scheduleAtFixedRate(
            new TimerTask(){
                public void run(){...} // periodic action 2
            }, 0, 500); // start now, period = 500ms
        ...
        periodic1.cancel(); // cancels when run() returns
        periodic2.cancel(); // cancels when run() returns
    }
}
```

Problems:

- No way to associate priorities with Timers
- SDK 1.3 documentation warns that Timer “does not offer real-time guarantees: it schedules tasks using the Object.wait(long) method”

What is a “priority inversion”?

- A situation in which a higher-priority thread is blocked/stalled involuntarily while a lower-priority thread is running

It is sometimes necessary (to enforce mutual exclusion)

- When the lower priority thread holds a lock that is needed by the higher priority thread

Scheduling policy affects the worst case blocking time

- Under some policies, a high priority thread may be blocked (stalled on a lock) during the execution of a lower-priority thread that does not hold the lock - *unbounded priority inversion*
- Priority Inheritance and Highest Lockers (Priority Ceiling) reduce this time considerably

```

class High extends Thread{
  private Object ref;
  High(Object ref){
    this.setPriority(10);
    this.ref=ref;
  }
  public void run(){ // periodic
    ... synchronized(ref){...} ...
  }
}

```

```

class Low extends Thread{
  private Object ref;
  Low(Object ref){
    this.setPriority(1);
    this.ref=ref;
  }
  public void run(){ // periodic
    ... synchronized(ref){...} ...
  }
}

```

```

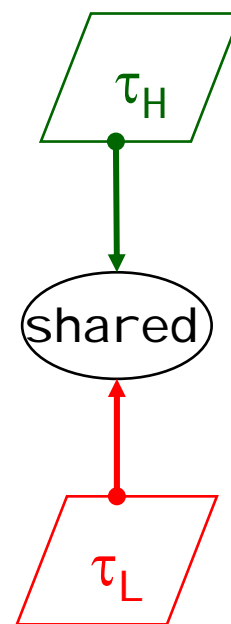
class Medium extends Thread{
  Medium(){ this.setPriority(5); }
  public void run(){ ... }
}

```

```

class PriorityInversion{
  public static void main(String[] args){
    Object shared = new Object();
    High  $\tau_H$  = new High(shared);
     $\tau_H$ .start();
    Low  $\tau_L$  = new Low(shared);
     $\tau_L$ .start();
    Medium  $\tau_M$  = new Medium();
     $\tau_M$ .start();
  }
}

```



Priority inversion example (2)



τ_H has high priority, τ_M medium priority, and τ_L low priority

- ① τ_L awakens and starts to run (the other two threads are blocked, waiting for the expiration of delays)
- ② τ_L starts to use a mutually-exclusive resource
 - Enters a monitor, locks a semaphore
- ③ τ_H awakens and preempts τ_L
- ④ τ_H tries to use the resource held by τ_L and is blocked, so τ_L resumes
 - This priority inversion is necessary
- ⑤ τ_M awakens and preempts τ_L
 - “Unbounded” priority inversion: τ_M indirectly prevents τ_H from running
- ⑥ τ_M completes, and τ_L resumes
- ⑦ τ_L releases the resource, is preempted by τ_H , which uses the resource
- ⑧ τ_H releases the resource
- ⑨ τ_H completes execution, allowing τ_L to resume
- ⑩ τ_L completes execution

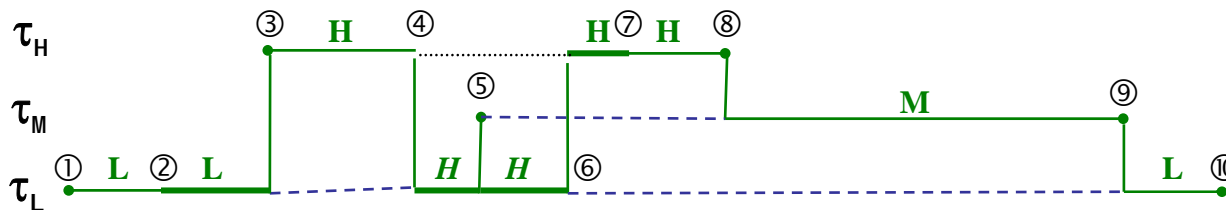
Priority Inheritance

- When a thread τ_H attempts to acquire a lock that is held by a lower-priority thread τ_L , τ_L *inherits* τ_H 's priority as long as it is holding the lock
 - Thus a medium priority thread τ_M cannot preempt τ_L while τ_L is holding the lock required by τ_H
- Applied transitively if τ_L is waiting for a lock held by a yet-lower-priority thread

Highest Lockers (Priority Ceiling Emulation)

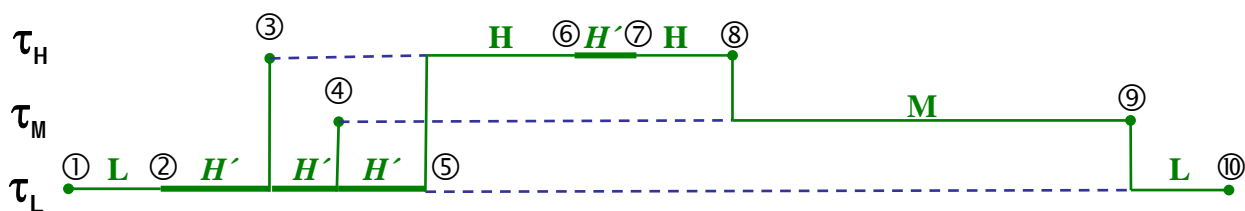
- Each lockable object has a *ceiling priority* equal to the highest priority of any thread that could lock the object
- When a thread acquires a lock, its priority is boosted to the object's ceiling
 - Run-time error if the thread's priority exceeds the object's ceiling
- When it releases the lock, its priority is reset
- Thus while holding a lock, a thread executes at a priority higher than or equal to that of any thread that needs the lock

Priority Inheritance example



- ① τ_L awakens and starts to run at priority L
- ② τ_L starts to use a mutually-exclusive resource
- ③ τ_H awakens, preempts τ_L and runs at priority H
- ④ τ_H tries to use the resource held by τ_L and is blocked, τ_L resumes
 - At this point τ_L inherits τ_H 's priority (H)
- ⑤ τ_M awakens *but does not preempt* τ_L
 - This avoids the unbounded priority inversion
- ⑥ τ_L releases the resource, reverts to its pre-inheritance priority L, and is preempted by τ_H , which can then use the resource
- ⑦ τ_H releases the resource
- ⑧ τ_H completes execution, allowing τ_M to execute
- ⑨ τ_M completes, allowing τ_L to resume
- ⑩ τ_L completes execution

Effect: a thread holding a lock executes at the max priority of all threads currently requiring that lock



- ① τ_L awakens and starts to run at priority L
- ② τ_L starts to use a mutually-exclusive resource with ceiling $H' \geq H$, and runs at priority H'
 - This will prevent unbounded priority inversion
- ③ τ_H awakens *but does not preempt* τ_L
- ④ τ_M awakens *but does not preempt* τ_L
- ⑤ τ_L releases the resource, reverts to its pre-ceiling priority L, and is preempted by τ_H which then runs at priority H
- ⑥ τ_H starts using the resource with ceiling $H' \geq H$, runs at priority H'
- ⑦ τ_H releases the resource and reverts to priority H
- ⑧ τ_H completes execution, allowing τ_M to execute
- ⑨ τ_M completes, allowing τ_L to resume
- ⑩ τ_L completes execution

Effect: a thread holding a lock executes at a priority no less than that of any thread that might need the lock

Advantages

- ☺ Supported by many RTOSes
- ☺ Only change priority when needed (thus no cost in common case when resource not in use)

Disadvantages

- ☹ Task may be blocked once for each shared resource that it needs (“chained blocking”)
- ☹ Implementation may be expensive
 - Task’s priority is changed by an action external to the task

Other properties

- Sacrifices responsiveness for predictability
 - A task may be prevented from running in order to guarantee that deadlines are met overall

Advantages

- ☺ If no thread can block while holding the lock on a given shared object, then a queue is not needed for that object
 - ☺ In effect, the processor is the lock
 - ☺ Prevents deadlock (on uniprocessor)
- ☺ Ensures that a task is blocked only once each period, by one lower priority task holding the lock

Disadvantages

- ☹ Fixed ceilings not appropriate for applications where priorities need to change dynamically
- ☹ Requires check and priority change at each call
 - Overhead even if object not locked
 - But this is inconsequential in the queueless case

Other properties

- If ceiling high, effect \cong disabling task switching
- Sacrifices responsiveness for predictability

There is a potential issue of storage leakage on assignment to a variable of a reference type

- May make inaccessible the object previously referenced

```
while(...){ Point p = new Point(0, 0);... }
```

In many languages this would be a problem

- Solutions involve unchecked storage reclamation, or provision of a finalization operation for the type/class

The JVM performs automatic reclamation (“garbage collection”)

- There is no language facility to explicitly deallocate an object
- The garbage collector operates (conceptually) as a background thread
 - It generally executes only when the program is otherwise suspended (e.g. waiting for user input)
 - It may also be invoked implicitly if there is not enough storage to fulfill an allocation request
- If you want to indicate that a reference to an object is no longer needed, then assign `null` to the variable
- You can also invoke the garbage collector directly

Main concepts

- At any point, the values that can be manipulated are the *roots* (stack, static variables) and the *live* heap objects
 - *Live* = directly or indirectly accessible from a root
- *Free space* = heap locations available for future allocations
- *Garbage* = objects that are neither live nor free

General issues related to dynamic memory management

- Storage leakage if garbage is generated and not freed
- Fragmentation if free space is not compacted
- Dangling reference if a live object is freed

Potential costs associated with garbage collection

- Time
 - Reclaiming objects, allocating objects
 - Assigning to a pointer, reading the value of a pointer
 - Delays imposed by the execution of the collector
- Space
 - Per heap object
 - Overhead intrinsic to collector algorithm

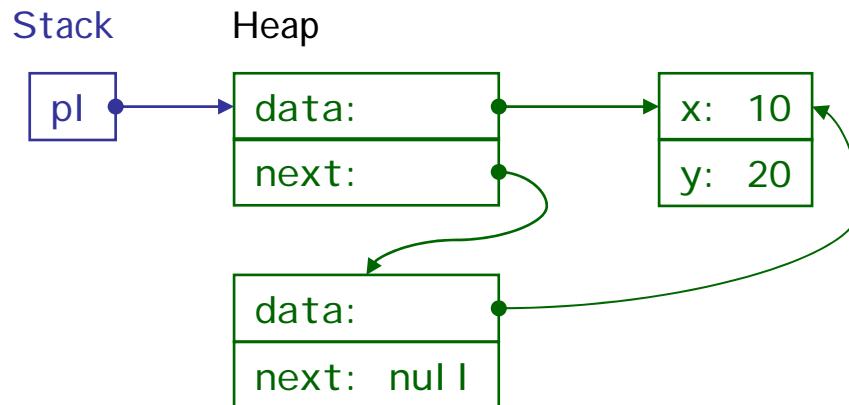
Example of live objects and garbage

```

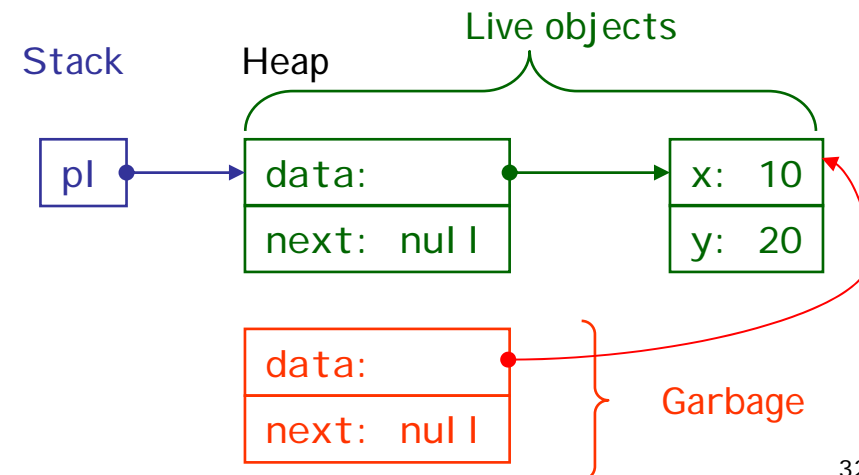
class PointList{
  Point data;
  PointList next;
  PointList(Point data, PointList next){
    this.data=data;
    this.next=next;
  }
public static void main(String[] args){
  PointList pl =
    new PointList(new Point(10, 20), null);
  pl.next =
    new PointList( pl.data, null); // (1)
  pl.next = null;                // (2)
}
}

```

State after (1)



State after (2)



Basics

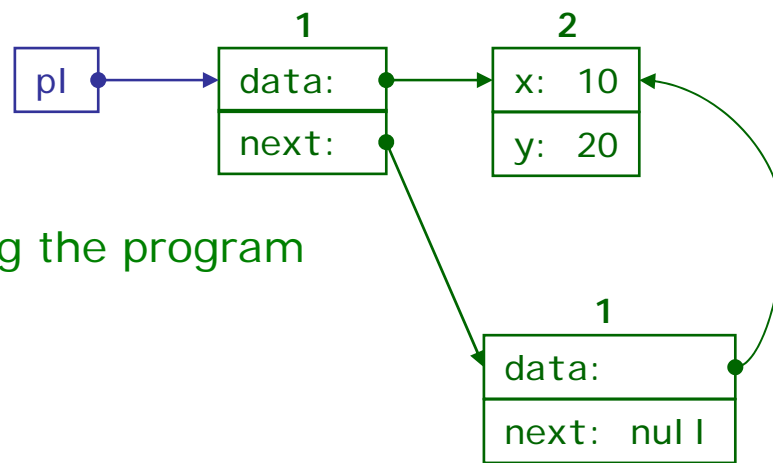
- Each object maintains a count of the number of pointers to that object, from roots or heap objects
- When a reference count is 0, the object can be freed

Advantages

- Conceptually simple
- Immediate reuse of short-lived objects
- Distributes overhead versus suspending the program

Disadvantages

- Does not detect cyclic structures
- Space penalty on each object
- Overhead on object allocation and pointer assignment to maintain the reference count
- Need special processing to avoid unbounded trace of an object whose count goes to 0
- Fragmentation



Basics

- *Mark* phase: global traversal of all live objects, starting from roots
- *Sweep* phase: any heap locations not marked are freed

Advantages

- Handles cycles
- No overhead on pointer manipulation
- No extra space per object

Disadvantages

- Computation halts during collection \Rightarrow not for real-time
- Fragmentation unless combine with compaction
- Sweep phase time proportional to size of heap
- Thrashing if heap occupancy high

Basics

- Heap comprises two semi-spaces, one containing current data and the other containing obsolete data
- Collector flips the two spaces, traverses the active objects in old space (*Fromspace*) and copies them into new space (*Toospace*), preserving object topology (sharing)

Advantages

- Free space is linear
 - Allocation cost low
 - Variable-sized objects do not present a problem
- Automatically compacts (avoids fragmentation)
- Cost proportional to live data versus whole heap (but see below)

Disadvantages

- Need to reserve twice the space actually needed
- Copying is more expensive than bit marking in Mark-Sweep
- Long-lived data objects are copied repeatedly
- Computation halts during collection ⇒ not for real-time

Basics

- “Weak generational hypothesis”: most objects die young
 - Focus reclamation on the objects created most recently
 - Collection triggered when number of bytes allocated since previous collection exceeds some threshold
- Segregate objects by age into heap regions (generations)
- Younger generation collected more frequently than older
- Young objects that survive some number of collections (typically 1 or 2) can be promoted to an older generation
- Most generational algorithms use copying, but mark-sweep possible

Advantages

- Pauses shorter since less data traced and copied
- Collection may be scheduled to minimize program disruption

Disadvantages

- Determining roots is more difficult, since pointers (especially from old to young) may be inter-generational
 - “Write barrier” overhead on pointer assignment
- “Tenured garbage” not reclaimed immediately

Goal is to guarantee *worst-case* pause time, for usability in real-time applications

- Generational technique tries to optimize *average* pause time

***Incremental* collector is interleaved with user program**

- Example: reference counting

***Concurrent* collector runs as a separate thread**

In both cases, collector must ensure heap state consistency

- The collector may temporarily treat some garbage as live (“floating garbage”)
- “Write barriers” typical with mark-sweep collectors
- “Read barriers” typical with copying collectors

User needs to ensure that garbage is not generated more quickly than it can be collected

In principle, an incremental or concurrent collector can be used in a real-time program

- But hardware support has been required for adequate performance

No Garbage Collection

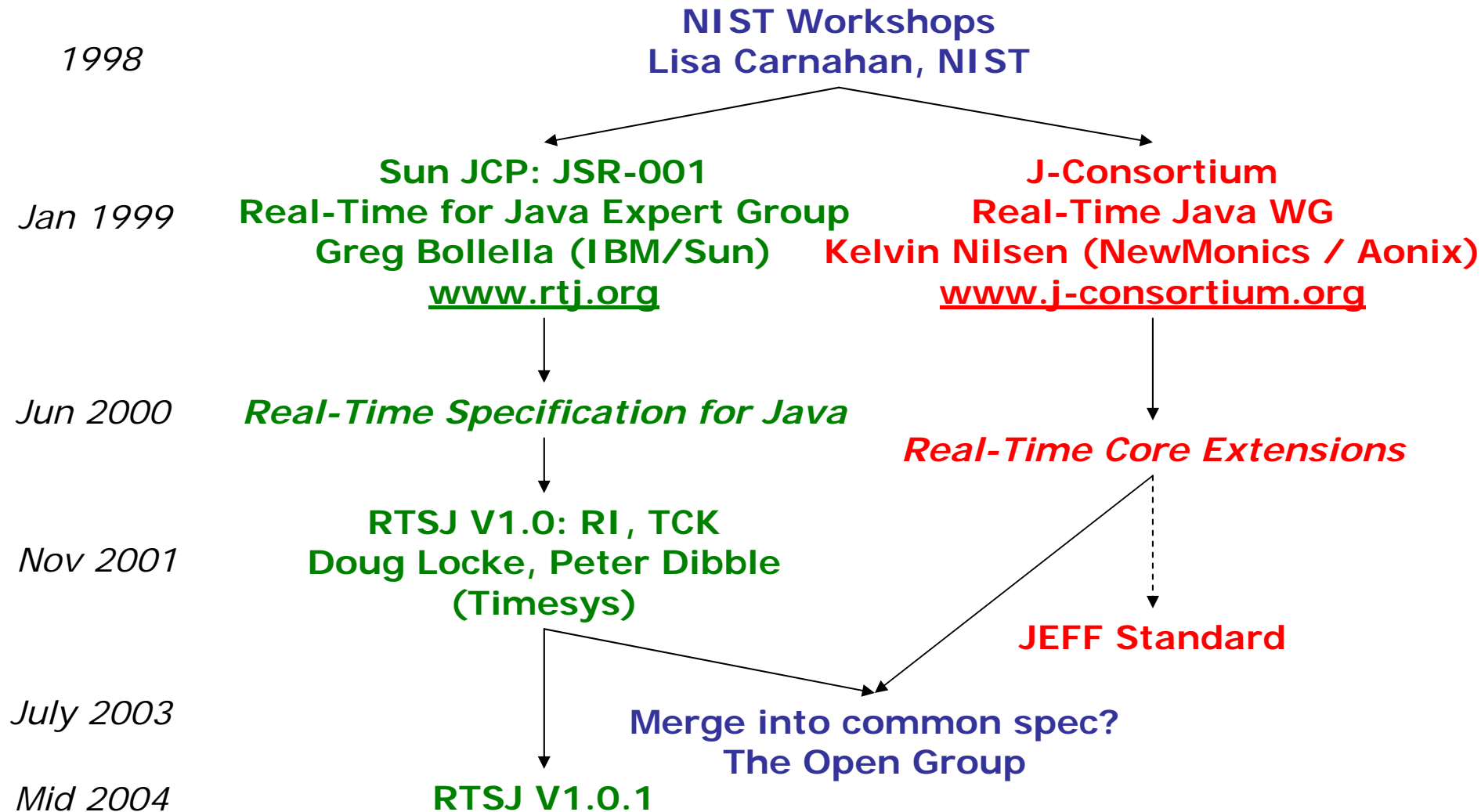
- Require that all allocations be performed at system initialization
- Common in many kinds of real-time applications
- Difficult in Java: all non-primitive data are dynamically allocated

Real-Time Garbage Collector

- Techniques exist that have predictable / bounded costs
- But programmer still needs to ensure that allocation rate does not exceed rate at which GC can reclaim space
- Also, in the absence of specialized hardware, such techniques tend to introduce high latencies
 - GC needs to run at high priority, since if it is preempted while the heap is in an inconsistent state, bad things can happen

Hybrid approach

- For low latency, allow a thread to preempt GC if the thread never references the heap
 - In general, need run-time check on each heap reference
- Allow a thread to allocate objects in a scope-associated area
 - Area flushed at end of scope/thread



Focus here will be on the Real-Time Specification for Java

Initial development (February 1999 – June 2000)

“Real-Time for Java Expert Group”, JSR-001

Greg Bollella (IBM / Sun)

James Gosling (Sun)

Ben Brosgol (Aonix / AdaCore)

David Hardin (Rockwell Collins / aJile)

Peter Dibble (Microware / Timesys) Mark Turnbull (Nortel Networks, Canada)

Steve Furr (QSSL, Canada)

Completion of V1.0 (July 2000 – November 2001)

Peter Haggart (IBM)

Doug Locke (Timesys)

Peter Dibble (Timesys)

Ben Brosgol (AdaCore)

Greg Bollella (Sun)

Rudy Belliardi (Schneider Automation)

Timesys Reference Implementation Team

RTJ-Technical Interpretation Committee (Nov. 2001- present)

Peter Dibble (Timesys)

Rudy Belliardi (Schneider Automation)

Greg Bollella (Sun)

Andy Wellings (Univ. of York, UK)

Ben Brosgol (AdaCore)

David Holmes (DLTeCH, Australia)

Applicability to particular Java environments

- Usage not restricted to particular versions of the JDK

Backward compatibility

- Existing Java code can run on RTSJ implementations

“Write Once, Run Anywhere”

- Important goal but difficult to achieve for real-time programs

Current practice vs. advanced features

- Address current real-time practice but allow implementations to include advanced features

Execution predictability

- Highest priority goal, may require compromise of performance/throughput

No syntactic extension

- No new keywords or other non-standard syntax

Allow variation in implementation decisions

- Recognize that implementations differ in terms of tradeoffs
- Do not mandate specific algorithms

Concurrency

- Class `RealTimeThread` extends `java.lang.Thread`
- Flexible scheduling framework, fixed-priority preemptive scheduler

Memory management

- “Immortal”, “Scoped” memory augment garbage-collected heap
- “NoHeap realtime thread” can preempt GC

Synchronization (priority inversion management)

- Priority inheritance (required), priority ceiling emulation (optional)

Asynchrony

- Asynchronous Event Handling
- Asynchronous Transfer of Control

Time and timers

- High-resolution time (absolute, relative)
- Timers (periodic, one-shot)

Low-level features

- Specialized kinds of “physical” memory
- “Peek/poke” of primitive data in “raw” memory



General concept of “schedulable object”

- Realtime thread or asynchronous event handler
- Arguments to constructor establish various “parameters”
 - Scheduling parameters (e.g. priority)
 - Release parameters (cost, deadline, overrun/miss handlers)
 - Periodic parameters (start, period)
 - Aperiodic parameters (min interarrival time for sporadic)

Initial default scheduler (“base scheduler”)

- Must support at least 28 distinct priority values, beyond Java’s 10
- Preemptive, fixed priority, FIFO within priority

Support for feasibility analysis (optional)

- Implementation can query release parameters to determine if a set of schedulable objects can satisfy some constraint

Flexibility

- Implementation can install arbitrary scheduling algorithms
- Users can replace these dynamically, can have different schedulers for different schedulable objects


```
public class MyPeriodic extends RealtimeThread{
    public MyPeriodic( int priority,
                      int periodMillis,
                      int costMillis,
                      int deadlineMillis ){
        super( new PriorityParameters( priority ), // SchedulingParameters
              new PeriodicParameters(           // ReleaseParameters
                null, // 1st release is at start
                new RelativeTime( periodMillis, 0 ),
                new RelativeTime( costMillis, 0 ),
                new RelativeTime( deadlineMillis, 0 ),
                null, // No cost overrun handler )
                null ) // No deadline miss handler
        );
    }
    public void run(){
        while (true) {
            ... // Work done during each release
            RealtimeThread.waitForNextPeriod(); // Block till next release
        }
    }
}
```

```
MyPeriodic mp = new MyPeriodic( 30, 100, 20, 20 );
mp.start(); // Triggers initial release for mp
```

Methods related to feasibility analysis

```
public boolean addToFeasibility()  
public boolean removeFromFeasibility()
```

- Include / exclude this schedulable object from the Scheduler's feasibility analysis

Methods for accessing "parameters" objects

- "Getter/setter" methods for SchedulingParameters, ReleaseParameters, MemoryParameters

Methods for accessing a scheduler

```
public Scheduler getScheduler()  
public void setScheduler(Scheduler scheduler)
```

- Get / set the Scheduler for this schedulable object

The run() method

- The executable logic for this schedulable object

Extends java.lang.Thread and implements Schedulabl e

Constructors

- General-purpose constructor

```
public Real timeThread( SchedulingParameters    scheduling,
                       ReleaseParameters       release,
                       MemoryParameters       memory,
                       MemoryArea             area,
                       ProcessingGroupParameters group,
                       java.lang.Runnable     logic)
```

- Constructors with default null values, interpreted based on default Scheduler at the time of the Real timeThread's construction

```
public Real timeThread()
public Real timeThread(SchedulingParameters scheduling)
public Real timeThread(SchedulingParameters scheduling,
                       ReleaseParameters  release)
```

Methods relevant to periodic RealtimeThreads

- Suspend until next period

```
public static boolean waitForNextPeriod()  
public static boolean waitForNextPeriodInterruptible()
```

- If not in deadline miss condition
 - Block until next release (start of next period)
 - Then return true
 - If in deadline miss condition
 - Block until released
 - Start of next period (if no handler), or when a handler calls `schedulePeriodic()`
 - Then return false
 - Interruptible version allows exceptional awakening from blocked state
- Begin / stop unblocking the target periodic Real timeThread

```
public void schedulePeriodic()  
public void deschedulePeriodic()
```

Relative and absolute sleep() method

```
public static void sleep( HighResolutionTime time )
```

- HighResolutionTime has subclasses AbsoluteTime and RelativeTime
- Another sleep method takes a Clock parameter
- Not needed for periodic Real timeThreads

Other methods

- Real timeThread analog to Thread.currentThread()

```
public static Real timeThread currentReal timeThread()
```

- Real timeThread overriding of Thread.interrupt()

```
public void interrupt()
```

- To be discussed below (Asynchronous Transfer of Control)

Cost enforcement (optional)

- If schedulable object *t* exceeds cost:
 - Handler, if any, is released (it may increase budget)
 - *t* is suspended until its next release event, at which time it is given a fresh budget and continues

ReleaseParameters is abstract superclass of several classes that reflect different kinds of release characteristics

- A ReleaseParameters instance is bound to a schedulable object at construction time, may be replaced or updated in place dynamically

PeriodicParameters

- Constructor

```

public PeriodicParameters( HighResolutionTime start
                           RelativeTime      period
                           RelativeTime      cost
                           RelativeTime      deadline
                           AsyncEventHandler overrider
                           AsyncEventHandler misshandler )

```

- Methods

- "Getter" and "setter" methods for attributes
- public boolean setIfFeasible(
 RelativeTime period, RelativeTime cost, RelativeTime deadline)

- Usage

- Schedulable object constructed with PeriodicParameters has run() method that loops on RealTimeThread.waitForNextPeriod
- Release is implicit, managed by implementation

AperiodicParameters

- Associate with schedulable objects whose releases are not periodic
 - Especially useful for AsyncEventHandlers
- Constructor
 - `public AperiodicParameters(RelativeTime cost
RelativeTime deadline
AsyncEventHandler overrunHandler
AsyncEventHandler missHandler)`
- Methods
 - `public void setInitialArrivalTimeQueueLength(int initial)`
 - Getters and setters for ArrivalTimeQueueOverflowBehavior
 - Throw exception, ignore, replace, or save
 - `public boolean setIfFeasible (RelativeTime cost,
RelativeTime deadline)`
- Usage
 - AsyncEventHandler constructed with an AperiodicParameters has `handleEvent()` method that responds to the event
 - Release occurs when the event is fired
 - In a future version of the RTSJ, `RealTimeThread` may have a `waitForNextRelease` method, relevant for AperiodicParameters

SporadicParameters (subclass of AperiodicParameters)

- Associate with aperiodic schedulable objects with a minimum interarrival time between releases
- Constructor
 - `public SporadicParameters(RelativeTime minimumInterarrival, RelativeTime cost, RelativeTime deadline, AsyncEventHandler overrunHandler, AsyncEventHandler missHandler)`
- Methods
 - `public void setMinimumInterarrival(RelativeTime minimum)`
 - Getters and setters for `MitViolationBehavior`
 - Throw exception, ignore, replace, or save
 - `public boolean setIfFeasible(RelativeTime cost, RelativeTime deadline)`
- Usage
 - Same style as for `AperiodicParameters`

Purpose

- Admission control for feasibility analysis
- Control over garbage collector pacing
- Control over storage usage by associated schedulable objects

Constructors

- `public MemoryParameters(`
 `long maxMemoryArea, // limit on # bytes in initial memory area`
 `long MaxImmortal, // limit on # bytes in immortal memory`
 `long allocationRate) // limit on # bytes / second`
- `public MemoryParameters(long maxMemoryArea,`
 `long MaxImmortal)`

Methods

- Getters and setters for `maxMemoryArea`, `maxImmortal`, `allocationRate`
- `public boolean setMaxMemoryAreaIfFeasible(long maximum)`
- `public boolean setMaxImmortalIfFeasible(long maximum)`
- `public boolean setAllocationRateIfFeasible(long allocationRate)`

Purpose

- Assign maximum execution time budget per “period”, collectively to a group of schedulable objects
- Allows inclusion of group of aperiodic schedulable objects in feasibility analysis, but may also apply to periodic and sporadic
- Budget enforcement not required if implementation cannot monitor execution time

Constructor

- `public ProcessingGroupParameters (HighResolutionTime start, RelativeTime period, RelativeTime cost, RelativeTime deadline, AsyncEventHandler overrunHandler, AsyncEventHandler missHandler)`

Methods

- Getters and setters for attributes (constructor parameters)
- `public boolean setIfFeasible (RelativeTime period, RelativeTime cost, RelativeTime deadline)`

Constructor

```
protected Scheduler()
```

Methods (partial list)

- Set/get default scheduler (initially an instance of PriorityScheduler)

```
public static void setDefaultScheduler(Scheduler s)
public static Scheduler getDefaultScheduler()
```

- Get name of the scheduling policy

```
public abstract String getPolicyName()
```

- Feasibility analysis related methods

```
protected abstract boolean addToFeasibility(Schedulable s)
protected abstract boolean removeFromFeasibility(Schedulable s)
public abstract boolean isFeasible()
public abstract boolean setIfFeasible( Schedulable s,
                                       ReleaseParameters r,
                                       MemoryParameters m )
```

Constructor

```
protected PriorityScheduler()
```

Methods (partial list)

- Static method to return the base scheduler

```
public static PriorityScheduler instance()
```

- Methods to get min, max, and normal priority

```
public int getMinPriority()  
public int getMaxPriority()  
public int getNormPriority()
```

- Method to trigger execution of a Schedulable object as an AsyncEventHandler

```
public void fireSchedulable (Schedulables)
```

- Static methods to get min, max, and normal priority for a Thread

```
public static int getMinPriority( Thread t )  
public static int getMaxPriority( Thread t )  
public static int getNormPriority( Thread t )
```

Represented by the object `PriorityScheduler.getInstance()`

Priority range

- Must support at least 28 distinct priority values, beyond the 10 for regular Java threads
- The 10 regular priority values have arbitrary mapping to native priorities
- In previous example, priority 30 should have been expressed more portably as `PriorityScheduler.getInstance().getMinPriority() + 20`

Preemptive, fixed priority, FIFO within priority

- Scheduler only changes a thread's priority for priority inversion avoidance
- Application may change a thread's priority dynamically
- Logical model is an array (indexed by priority) of ready queues
- "Run till blocked or preempted"
 - "Timeslice within highest priority" not permitted for base scheduler but may be supplied by implementation as additional scheduler

Several actions cause a thread* to go to the tail of the ready queue for its active priority

- A blocked thread becoming ready, via `notify()` / `notifyAll()`
- A sleeping thread becoming ready after delay expires
 - Including `waitForNextPeriod`
- A thread whose priority is changed via `PriorityParameters.setPriority()`
- A thread that performs a `yield()`

The RTSJ does not specify placement in the ready queue for a thread that is preempted or suffers loss of inherited priority

- Advises placement at head of queue
- Implementation must document its algorithm
- Non-determinism motivated by concerns over what RTOSes do in practice

** The term “thread” here includes schedulable objects*

Advantages

- General, flexible, extensible framework
 - Various kinds of release characteristics reflecting common uses
 - “Hooks” for implementation-defined schedulers
 - Provision for different schedulers governing different threads
 - Allowance of dynamic replacement of schedulers
- On-line feasibility analysis with execution-time budget enforcement

Disadvantages

- Clumsy syntax
 - Constructors with positional notation and no default values
 - Retrieving min, max, norm priority values for base scheduler
- Design shows signs of “rush to publish”, and design bugs are difficult to correct without introducing incompatibilities
 - Missing some useful elements
 - Class hierarchy has some awkwardness
- Implementation-defined aspects of base scheduler
- Absence of other useful scheduling policies (round robin, EDF)
- Possible performance issues

Issue

- Garbage Collector may introduce delays that interfere with predictability and/or performance

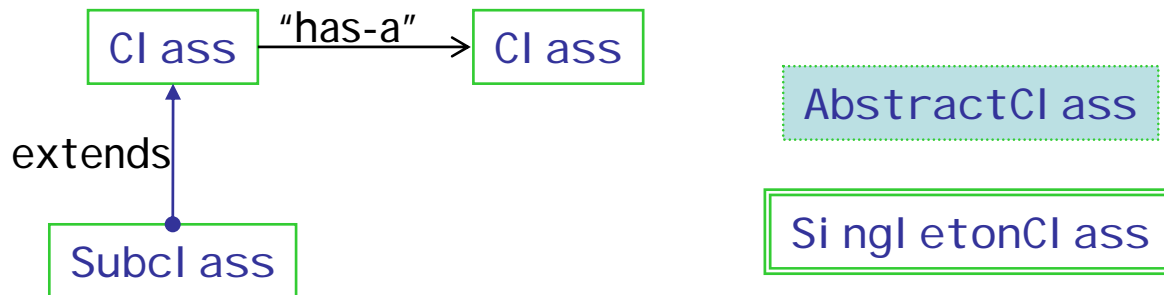
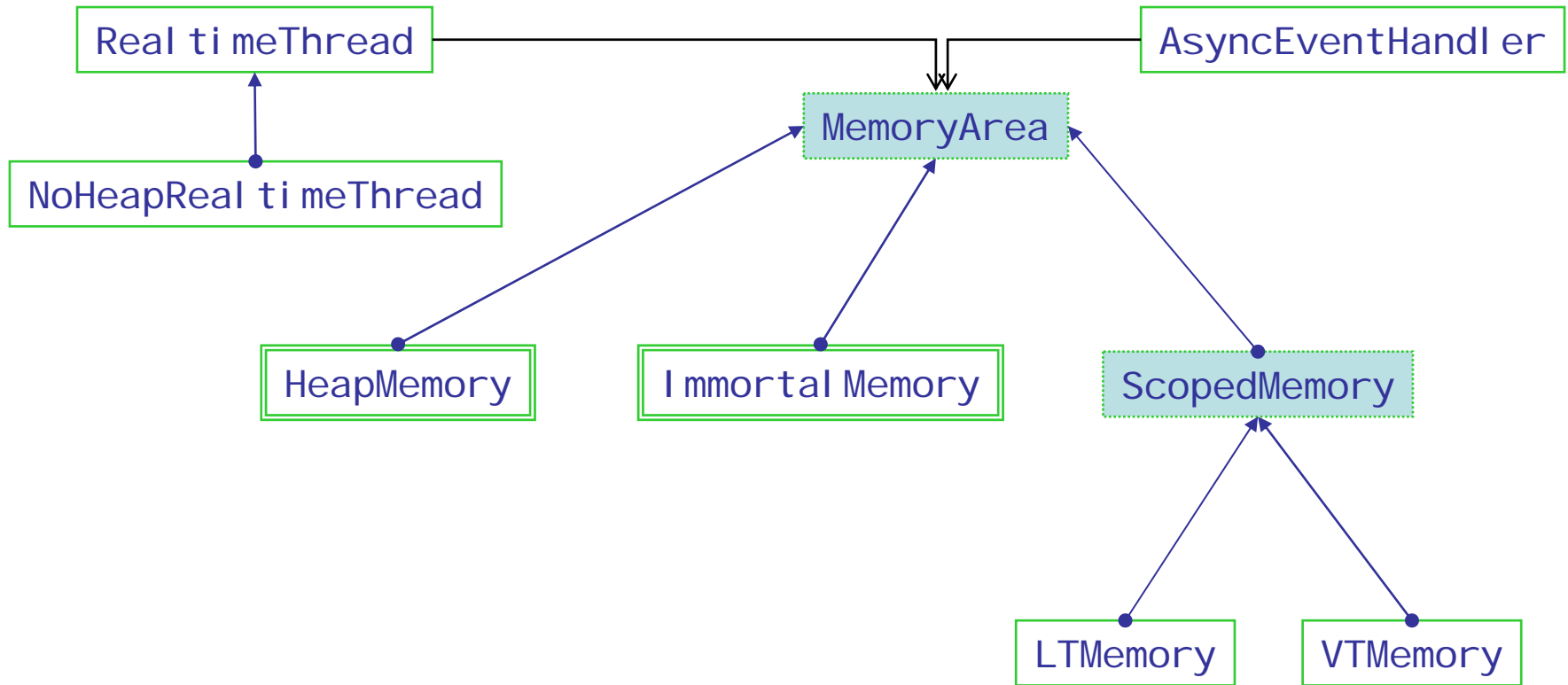
Goals

- Provide an efficient, general, predictable approach to memory management that can supplement the regular GC
- Do not compromise Java safety (i.e., no explicit “free”)

Approach

- General concept of a “memory area” (context for object allocation)
 - Garbage-collected heap (singleton)
 - Non-GC’ed areas
 - Immortal memory (singleton)
 - Scoped memory (user-instantiable)
- Subclass of real-time threads that never reference the heap
 - They may preempt the GC and, in general, do not suffer GC latency
 - They may reference immortal or scoped memory

Memory management classes (partial list)



Definition

- The *closure* (of a method or constructor $f()$ invoked by a thread or schedulable object t) is the sequence of statements executed by t for $f()$, including those in all methods or constructors called directly or indirectly from f , until f returns
 - This does not include the statements in the `run()` method of threads / schedulable objects that are started from $f()$, since such statements are not executed by t

You can bind the closure of the `run()` method for a `Runnable` object r to a memory area ma

- Within a given thread or schedulable object: `ma.enter(r)`
- Across multiple schedulable objects: pass ma as constructor parameter
- All object allocations during the closure will use ma , unless overridden by an inner binding to a different memory area

You can also bind a single constructor invocation to a memory area ma

- `ma.newInstance(class)`
- `ma newArray(class, length)`

Regular Java heap

- Singleton class `HeapMemory`
- Single instance given by static method `HeapMemory.getInstance()`
- Subject to the implementation's Garbage Collection algorithm

"Immortal Memory"

- Singleton class `ImmortalMemory`
- Single instance given by static method `ImmortalMemory.getInstance()`
- Objects are not relocated or deallocated
 - Released only at program termination
 - Whether finalizers are run is implementation dependent
- May point to heap and vice versa
 - Thus GC needs to be able to scan immortal memory

Motivation

- In Java, only primitive data go on the stack, but most objects allocated in a method are only used locally in that method
 - Arrays and class instances thus need to be garbage collected
- Allow user to specify that some objects are “stackable”?
 - Does not work for objects that are returned as method results
- “Escape analysis” to automatically store objects on stack when safe?
 - Useful as an optimization, but not sufficient since some objects might still need to go on heap

Approach

- Allow user to construct non-GC'ed reusable memory areas that will be used for allocations during execution of specific code regions (closures)
- Prevent dangling references

Classes

- Abstract class `ScopedMemory`
- Subclass `LTMemory` requires implementation to allocate / initialize with a linear time bound (based on object size)
- Subclass `VTMemory` allows “varying” (i.e. non-linear) time bound

Basic usage similar to heap and immortal

- May be shared among / bound to several schedulable objects by being passed to their constructors
 - Used for allocations during the closure of the threads' `run()` methods
- May be used locally by a given schedulable object
 - Invoking `sm.enter(r)` activates the scoped memory `sm` for Runnable `r`'s `run` method and synchronously invokes `r.run()`
 - Used for allocations during the closure of `r.run()`

Distinctive features of scoped memory

- User creates scoped memory areas with specified sizes
- Each scoped memory area has an associated “backing memory” in which the allocated objects are stored
- A reference count scheme (number of referencing schedulable objects) establishes when a scoped memory area can be reset
 - Objects finalized, area emptied when last schedulable object bound to the area terminates, or when outermost `enter` returns
- Reference assignment is restricted, to prevent dangling references

```
public abstract class MemoryArea{
    protected MemoryArea( long sizeInBytes );

    protected MemoryArea( long sizeInBytes, Runnable logic );
    // logic.run() is invoked whenever enter() is called on the constructee

    public void enter( Runnable r );
    // Binds this to the execution of r.run(), as the area used for allocations

    public synchronized Object newInstance(Class c)
        throws IllegalAccessException, InstantiationException, OutOfMemoryError;

    public synchronized Object newArray(Class c, int n)
        throws IllegalAccessException, InstantiationException, OutOfMemoryError;

    public static MemoryArea getMemoryArea(Object o);

    public long size();
    // current size in bytes (including unused space)

    public long memoryConsumed();
    // size (in bytes) of memory allocated to objects

    public long memoryRemaining();
    // approximate size (in bytes) of memory available for future allocation
}
```

For a single allocation

```
ma.newInstance(someClass) or ma.newArray(someClass, n)
```

For a section of code, including methods called from there

```
class MyRunnable implements Runnable{  
    public void run(){ ... }  
}  
MyRunnable r = new MyRunnable();  
ma.enter(r); // invokes r.run() synchronously
```

Alternative style with anonymous class

```
Runnable r = new Runnable(){  
    public void run(){  
        ... // code whose "new" objects are allocated in ma  
    }  
}  
ma.enter(r);
```

Another style using an anonymous object from an anonymous class

```
ma.enter(new Runnable(){  
    public void run(){  
        ... // code whose "new" objects are allocated in ma  
    });
```

```
class Sensor extends RealTimeThread{

    final Position ps;

    final LTMemory ltm = new LTMemory(10*1024, 10*1024); // Initial size = max size = 10KB

    Sensor(Position p){
        super(new PriorityParameters(...), new PeriodicParameters(...));
        ps = p;
    }
    public void run(){
        double x, y;
        while (true){
            ltm.enter(new Runnable(){
                public void run(){
                    double x = InputPort.read(1);
                    double y = InputPort.read(2);
                }
            }); // Any objects allocated during read() are finalized, and the area is cleared
            synchronized(ps){ ps.x=x; ps.y=y; }
            try{
                this.waitForNextPeriodInterruptible();
            }
            catch( InterruptedException e){
                return;
            }
        }
    }
}
```

From To	Reference to Heap	Reference to Immortal	Reference to Scoped
Heap	Yes	Yes	No
Immortal	Yes	Yes	No
Scoped	Yes	Yes	See (1)
Local Variable	Yes	Yes	See (2)

(1) Yes if object containing the destination field is in same or inner (more recently entered) scoped memory

(2) Yes if destination variable is local to block in same or inner scope

Implementation may detect violations at compile time rather than throwing an exception at run time

In many potential application domains, such compile-time analysis will be important if the RTSJ is to be competitive with existing technologies

```
ma = new LMemory(i ni tSi ze, maxSi ze);
```

- Allocates an LMemory instance in the current memory area
- Creates ma's "backing memory" based on the size parameters
 - "Backing memory" is outside the Java realm (e.g. mal l oc'ed)

The reference count for the memory area instance is incremented when:

- ma is entered: `ma.enter(r)`
- A schedulable object with ma as memory area parameter is started

The reference count is decremented when:

- Control returns from `ma.enter(r)` (i.e., control returns from `r.run()`)
- A schedulable object which has ma as memory area parameter terminates

When the reference count goes from 1 to 0:

- Finalizers are run for all objects in the memory area
- The memory area is "emptied" (i.e., the entire backing memory is again made available)

Note: the backing memory is not freed when the reference count goes to 0

Example: a shared ScopedMemory (1)

```
class Sensor extends RealTimeThread{
    final Position ps;
    Sensor(Position p, MemoryArea ma){
        super(new PriorityParameters(...),
              new PeriodicParameters(...),
              null, null, ma, null, this);
        ps = p;
    }

    public void run(){
        // All news use the "backing memory" associated with ma
        while (true){
            double x = InputPort.read(1);
            double y = InputPort.read(2);
            synchronized(ps){ ps.x=x; ps.y=y; }
            try{
                this.waitForNextPeriodInterruptible();
            }
            catch (InterruptedException e){
                return;
            }
        }
    }
}

class Reporter extends RealTimeThread{... }
// analogous
```

```
class Test{
  public static void main(String[] args){
    Position p = new Position();
    final LTMemory ltm =
      new LTMemory(10*1024, 10*1024);
    // Initial size = max size = 10KB

    Sensor s = new Sensor(p, ltm);
    Reporter r = new Reporter(p, ltm);
    s.start(); r.start();
    ...
    s.interrupt();
    // terminate s at start of next period
    r.interrupt();
    // terminate r at start of next period
  }

  // When s, r, and the main thread have terminated,
  // any objects allocated in ltm's backing memory
  // are finalized, and the area is cleared
}
```

Finalization / refresh of backing memory when a scoped memory area's reference count goes to 0 does not deallocate the backing memory

```
void foo(){
    LTMemory m = new LTMemory( 100000, 100000 );
    m.enter(
        new Runnable(){
            public void run(){...}
        }
    );
    // Finalize all objects in m's backing memory, and reset to empty
}
while ( someCondition() ){
    foo();
}
```

- Each time foo is invoked, a new MemoryArea instance and 100KB backing memory space is allocated, and the area is finalized / refreshed when enter returns

In general, the backing memory can only be freed when its associated instance has been allocated in an outer scoped memory area, as part of the finalization of this instance

The following style allows the same scoped memory area to be used repeatedly, with the area finalized / refreshed after each use

```
void foo(){
    LTMemory m = new LTMemory( 100000, 100000 );
    while ( someCondition() ){
        m.enter(
            new Runnable(){
                public void run(){...}
            }
        );
        // Finalize all objects in m's backing memory,
        // and reset the backing memory to empty
    }
}
```

Memory areas and schedulable objects

- The initial memory area for a schedulable object is defined either by a constructor parameter or by the creating thread's current memory area
- The memory areas currently active for a schedulable object form a stack

Restrictions on scoped memory allow “reference count” implementation

- A scoped memory area may be shared or entered, but not both
- No nested enters on the same scoped memory area

Regular Java threads may only use heap and immortal memory

- Constructing a Thread in a scoped memory area throws exception

Exceptions are thrown as usual

- Allocated in current memory area
- If attempt to propagate out of a scope, throw a `ThrowBoundaryError` exception (itself allocated in surrounding memory area)

Advantages

- Scoped memory provides general, predictable, non-GC'ed storage management
- SizeEstimator class gives implementation-independent mechanism for defining size of a memory area
- Safe (no explicit "free" / "unchecked deallocation")

Disadvantages

- Lots of objects need to be stored in immortal memory
 - Class objects, objects created during static initializers
- Reclamation of "backing memory" is implementation dependent
 - A method that allocates a scoped memory area may result in storage leakage if it is invoked in a loop
- Users need to be aware of memory management issues
 - Calling a library method when a scoped memory area is active may cause an assignment error (e.g., allocating an object and storing a reference in a static variable)
- Performance penalty for assignment checks, unless static analysis
- An exception propagating out of scoped memory areas will cause a cascade of `ThrowBoundaryError` exceptions

Monitor control policy allows user to select which policy governs which objects

- Semantics defined for default scheduler
- Distinction between active and base priority

Priority Inheritance is default policy

- May be changed by user at system startup

Priority Ceiling Emulation is also defined (but is optional)

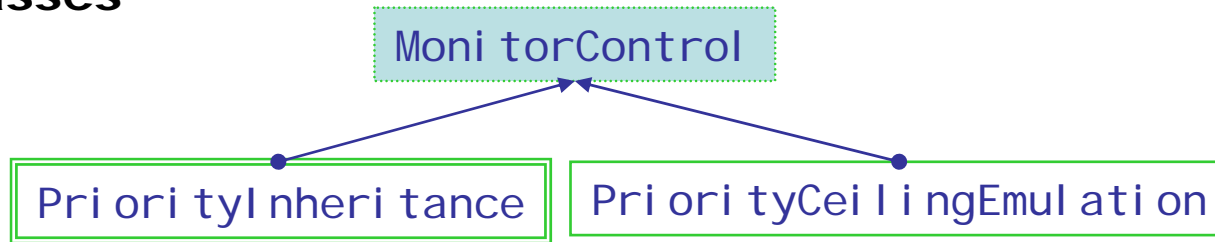
- Locking thread's priority is boosted to ceiling when lock acquired, reset when lock released
- Ceiling violation exception thrown if locking thread has higher priority than the ceiling
- No requirement for non-blocking as in Ada

Dynamic mechanism

- May change an object's monitor control policy, modify ceiling

“Wait-free queues” allow communication between a NoHeap Realtime Thread and a heap-using thread

Classes

Abstract class `MonitorControl`

- Methods to set policy globally or per-object

```

public static void setMonitorControl (MonitorControl policy)
public static void setMonitorControl (Object obj ,
                                     MonitorControl policy)
  
```

- Methods to get policy globally or per-object

```

public static MonitorControl getMonitorControl ()
public static MonitorControl getMonitorControl (Object obj )
  
```

Class `PriorityInheritance`

- Method to return the singleton instance

```

public static PriorityInheritance instance()
  
```

Class PriorityCeilingEmulation

- Uses
 - Throw CeilingViolationException if the synchronizing thread or schedulable object has a priority higher than the target object's ceiling
 - Immutable class, so need to setMonitorControl(obj) to change obj's ceiling
- Method to create policy object / return existing object with given ceiling

```
public static PriorityCeilingEmulation instance( int ceiling )
```

- More efficient than constructor
- Method to return ceiling for this policy object

```
public int getCeiling()
```

- Method to return max ceiling supported by base scheduler

```
public static int getMaxCeiling()
```

```
class Position{
  private double x, y;
  Position(double x, double y, int ceiling){
    this.x=x; this.y=y;
    MonitorControl.setMonitorControl(this, PriorityCeilingEmulation.instance(ceiling));
  }
  synchronized void setXY(double x, double y){
    this.x=x; this.y=y;
  }
  synchronized double[] getXY() {
    return new double[2]{x, y};
  }
}

class Sensor extends RealtimeThread{
  Sensor(Position p, int priority){...}
  ...
}

class Reporter extends RealtimeThread{
  Reporter(Position p, int priority){...}
  ...
}

class Test{
  public static void main(String[] args){
    Position p = new Position(0.0, 0.0, 16);
    Sensor s = new Sensor(p, 15);
    Reporter r = new Reporter(p, 15);
    s.start(); r.start();
    ...
  }
}
```

Each thread or schedulable object t has a *base priority* and an *active priority*

- Base priority is initially its priority at construction
- Updated immediately by several priority-setting methods

If t holds no locks, its base priority = active priority

- Priority-setting method moves t to tail of queue at new priority

If t holds one or more locks it has a set of *priority sources*

- *Active priority* is the max of the priorities from all sources

Source	Priority
t itself	Base priority for t
Each PCE-governed object locked by t	Max ceiling of all PCE locks
Each thread trying to synchronize on a PI-governed object locked by t	Max active priority of all such threads

- If the addition of a priority source increases t 's active priority, t is placed at the tail of the relevant queue at the new active priority
- If the removal of a priority source decreases t 's active priority, t should be placed at the head of the relevant queue at the new active priority (implementation advice, not required)

RTSJ does not prohibit synchronized code from blocking

- Thus queues / mutexes are needed to implement lock management

A more restricted form of PCE was considered for inclusion

- Throws an exception if a thread blocks while holding a lock

Advantages of non-blocking PCE

- More efficient implementation (no need for queues)
- Good model for interrupt handlers
- Interoperability with Ada protected objects

For simplicity, non-blocking PCE was omitted

- For example, would need to prevent assigning a non-blocking PCE lock to an object whose synchronized code can block
 - Otherwise exception will propagate out of synchronized code with the object in an inconsistent state

Implementation can add the policy but:

- Still needs to implement queue-based PCE
- Needs to ensure that the policy can only be assigned to objects that permit it

Problem with priority inheritance if a NoHeapReal timeThread t1 and a heap-using thread or schedulable object t2 attempt to synchronize on an object obj in ImmortalMemory or ScopedMemory

- The heap-using entity t2 locks obj
- The Garbage Collector preempts t2
 - May be triggered by “heap getting low” event
- The NoHeapReal timeThread t1 preempts the GC and tries to lock obj
- The heap-using entity t2 inherits the NoHeapReal timeThread's priority p1 (higher than the GC) and runs
- If t2 allocates a heap object the GC's data structures will be corrupted

Possible approaches to avoiding this problem?

- Throw an exception in t1 when it attempts to lock obj (no priority inheritance by t2)
 - Hard to program if a NoHeapReal timeThread can have an exception thrown just because it is synchronizing on a PI-governed object
- Defer the priority inheritance by t2 until the GC is in a safe state
 - May cause unbounded priority inversion for t1, by lower-priority NoHeapReal timeThreads whose priority exceeds that of the GC

Better solution involves a conceptual lock on the heap

- Accesses to the heap (by a heap-using thread or schedulable object, or by the GC) are controlled by a Priority Inheritance-governed lock
- In the above example, the GC owns the lock when it is preempted by the NoHeapReal timeThread t1
- The heap-using entity t2 inherits t1's priority p1 and runs until it attempts to use the heap
- When t2 tries to acquire the heap lock, it will be blocked and the GC will inherit t1's priority

Consequences of this solution

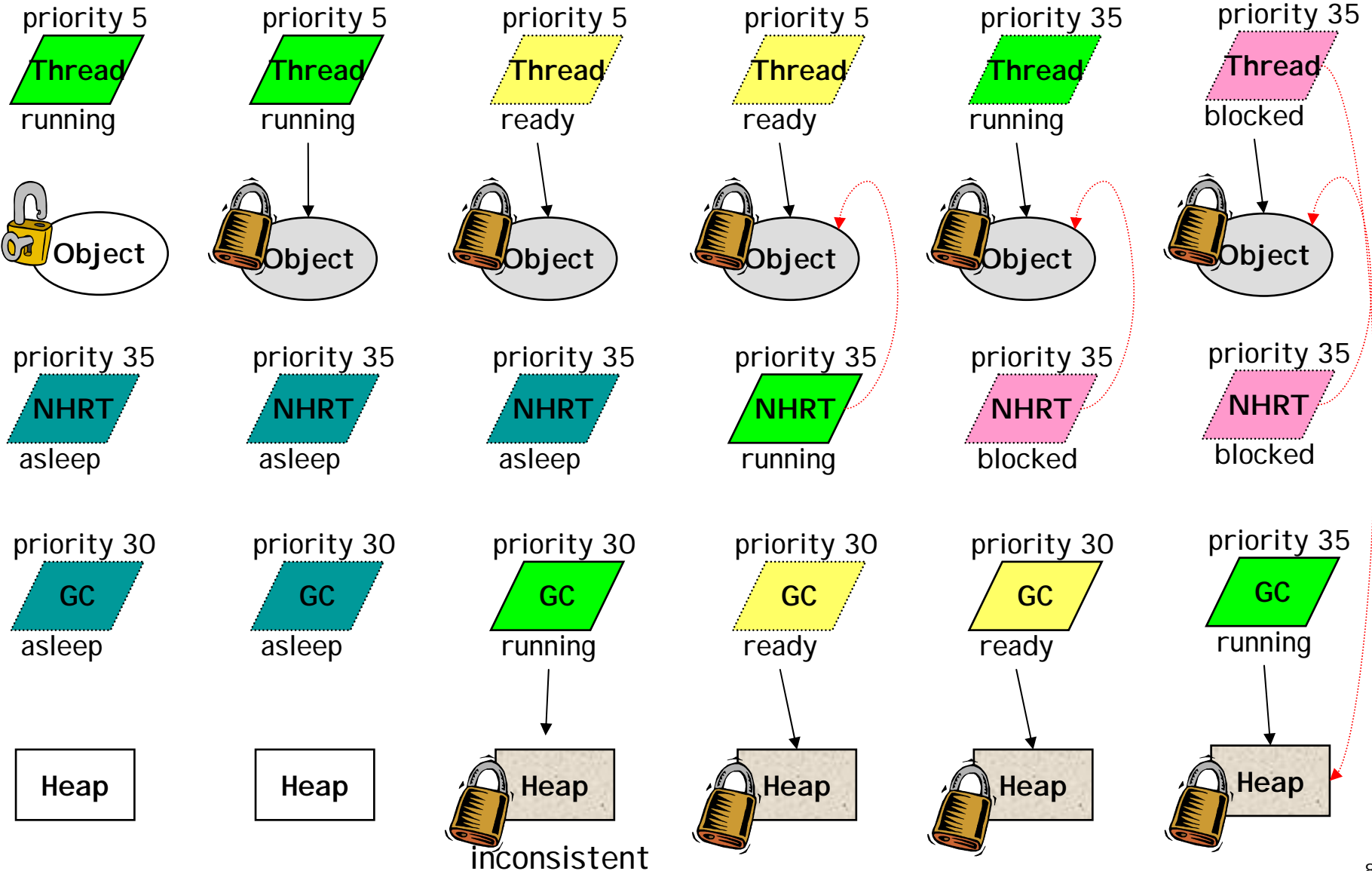
- Allows immediate priority inheritance, thus avoids unbounded priority inversions
- But it does introduce potential GC latency for a NoHeapReal timeThread

Alternative programming solution

- If the GC latency is unacceptable, use a "Wait-Free queue" (versus synchronizing on an object) for communication between a NoHeap-Real timeThread and a heap-using thread or schedulable object
 - Ensures that the NoHeapReal timeThread does not block

Communication between no-heap, heap threads

GC may introduce latency for a NoHeapRealtimeThread



Issue: how to share an object between a NoHeapReal timeThread and a heap-using thread or schedulable object

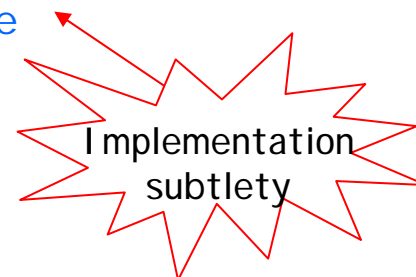
- Using synchronized code may lead to GC-induced priority inversions

Approach: two “wait-free” queues allow 1-writer multiple-reader or 1-reader multiple-writer communication

- Communication through bounded buffer of Objects, managed FIFO

WaitFreeWriteQueue class

- Single writer (generally a NoHeapReal timeThread)
- One or more readers (heap-using thread or schedulable object)
- `q.write(obj)` returns boolean
 - Non-synchronized, non-blocking
 - if `q.nonfull` then insert `obj` and return `true`
 - if `q.full` then return `false` (do not overwrite existing object)
 - In either event, awaken a waiting reader if there is one
- `q.read()` returns `Object`
 - Synchronized, potentially blocking
 - If `q.nonempty`, return object least-recently written
 - If `q.empty`, block (wakeup will be from eventual writer)



WaitFreeReadQueue class

- Single reader (generally a NoHeapRealTimeThread)
- One or more writers (heap-using thread or schedulable object)
- Control if reader blocks on queue empty, unblocks on non-empty
- `q.write(obj)`
 - Synchronized, potentially blocking
 - if `q` nonfull then insert `obj`
 - if `q` full then block until there is room
 - Awaken blocked reader (if any) if `q` appropriately constructed
- `q.read()` returns Object
 - Non-synchronized, non-blocking
 - If `q` nonempty, return object least-recently written
 - If `q` empty, return null
- `q.waitForData()`
 - Non-synchronized, potentially blocking
 - If `q` nonempty, return immediately
 - If `q` empty, block (wakeup will be from eventual writer)
 - Style: `q.waitForData(); obj = q.read();`



Advantages

- Highly flexible, with monitor control policies assignable per-object or globally
- Supports both priority inheritance and priority ceiling emulation
- Extensible (e.g. implementation may add new policy classes)
- Efficient approach to PCE policy instances (one per ceiling value)

Disadvantages

- No easy way to obtain optimization of non-blocking PCE
- Mixture of PI and PCE results in some semantic complexity and likely also some run-time overhead
- Calculating appropriate ceiling for an object can be tricky
- Wait-free queues may lead to storage leakage of immortal memory if not used carefully
 - User needs global knowledge of synchronization behavior
- Deprecated WaitFreeDequeue class shows signs of rushed design
- Optionality of PriorityCeilingEmulation interferes with portability

An *asynchronous event* for a thread t affects t without t 's direct request, at an unpredictable point in t 's execution

- May be generated externally (e.g. from an interrupt or an OS signal) or internally (from another thread)
- A thread registers one or more AsyncEventHandl er objects with an AsyncEvent object,
 - Scheduling parameters are supplied
- The handler is run (conceptually in an implementation-provided thread) in response to occurrence of an internal event or external "happening"

Handler can make its effect known to t in several ways

- Setting a variable polled by t
 - Example: input available
- Causing an asynchronous transfer of control in t
 - Example: time out and abandon computation

Goal

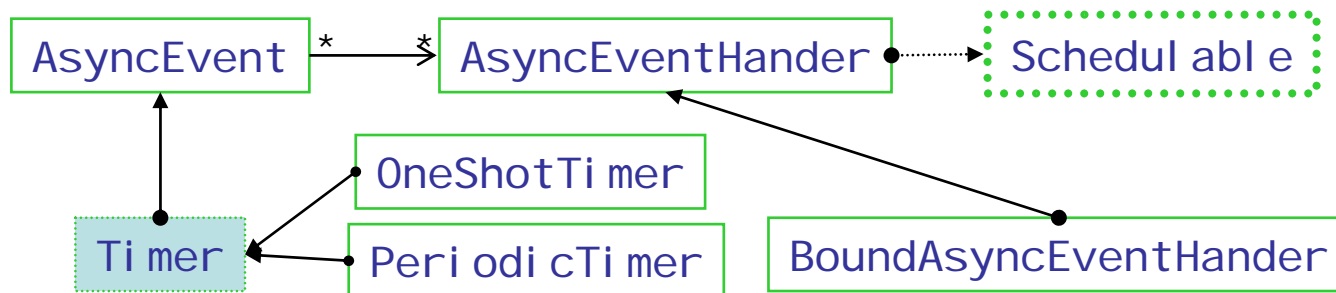
- Scale up to deal with large number of events and event handlers (tens of thousands), but with a smaller number of event handlers simultaneously active

Asynchronous Event Handler (“AEH”)

- Use for hardware interrupts or software “happenings”
- An AEH is a schedulable object but need not have a dedicated thread
- Override a method to implement the relevant event handling
- Associate one or more Asynch Event Handlers with an Asynch Event
 - Firing an AE → schedule associated AEHs
- Traditional Java “listener registration” model

Asynchronous Transfer of Control (“ATC”)

- Use for timing out on a computation, aborting a thread
- Methodologically questionable, and complicated to implement
 - Conflict between desire for ATC to be immediate, and the need for certain code to execute completely
- Extends `t.interrupt()` to real-time threads, throwing an exception not only when `t` is blocked but also when `t` is executing *asynchronously interruptible* (“AI”) code
 - Synchronized code, and methods lacking a special throws clause, are not AI



The AsyncEvent class

- Public no-arg constructor
- Methods (partial list)
 - public void addHandler(AsyncEventHandler handler)
 - public void removeHandler(AsyncEventHandler handler)
 - public void setHandler(AsyncEventHandler handler)
 - public void bindTo(String happening)
 - public void unbindTo(String happening)
 - public void fire()

The AsyncEventHandler class

- Constructors (partial set)

- ```
public AsyncEventHandler(SchedulingParameters scheduling,
 ReleaseParameters release,
 MemoryParameters memory,
 MemoryArea area,
 boolean nonheap,
 Runnable logic)
```
  - ```
public AsyncEventHandler( SchedulingParameters scheduling,
                          ReleaseParameters   release,
                          MemoryParameters   memory,
                          MemoryArea         area,
                          boolean            nonheap )
```

- Methods (partial list)

- public void handleAsyncEvent()
 - protected int getPendingFireCount()
 - protected int getAndClearPendingFireCount()
 - protected int getAndDecrementPendingFireCount()
 - Getters and setters for “parameters” objects
 - Methods related to admission control / feasibility analysis
- } called from handleAsyncEvent

AsyncEventHandlers are “registered” with an AsyncEvent by invoking **ae.addHandler(aeh)**

- Many-many relationship between events and handlers

An AsyncEvent ae is triggered by an invocation **ae.fire()**

- From a hardware interrupt handler or from a software thread
- No data passed
- Handler has event count in case of bursts

An AsyncEventHandler behaves conceptually like a thread (and thus has scheduling parameters etc)

- A BoundAsyncEventHandler has a dedicated thread

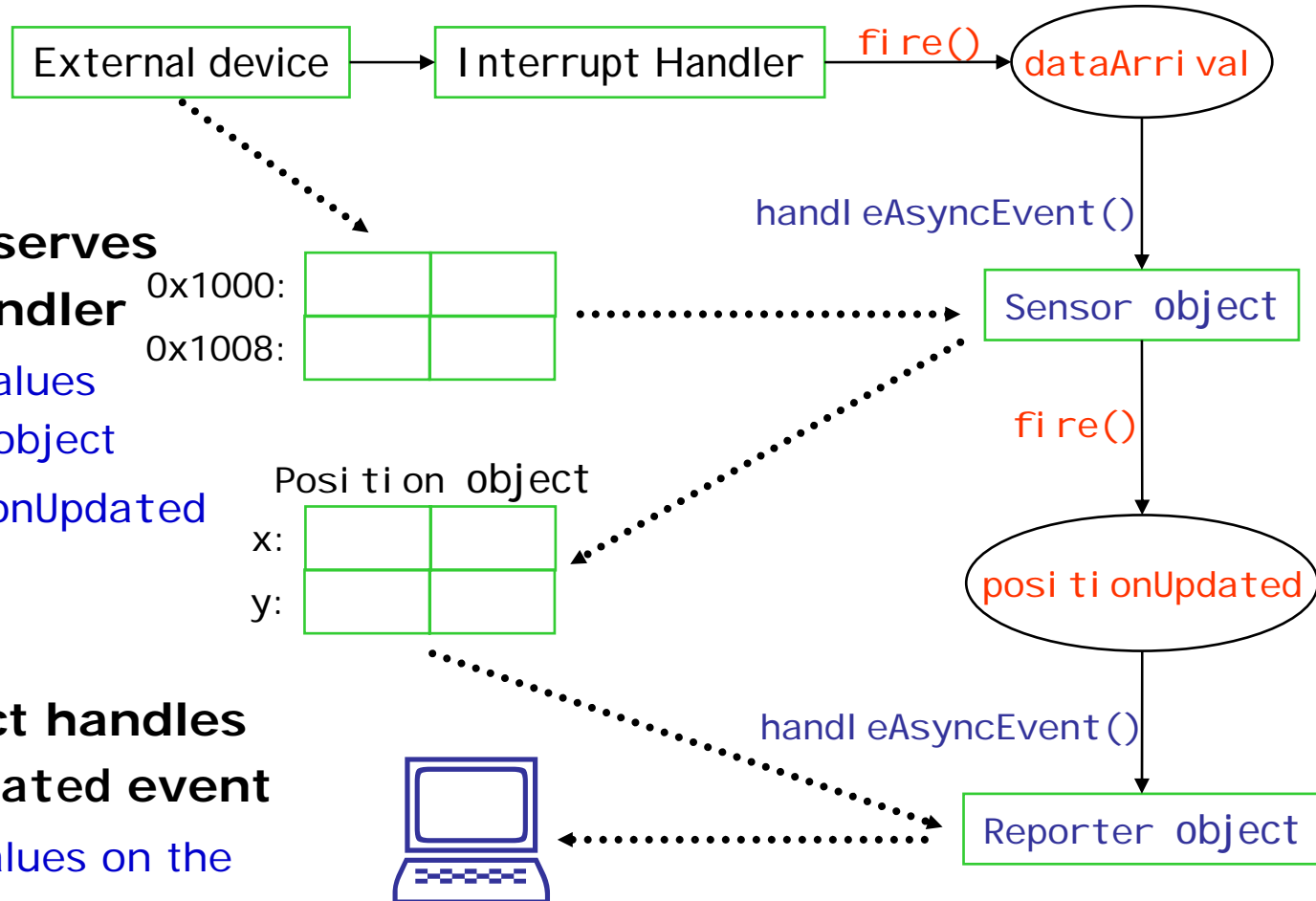
When an AsyncEvent ae is fired, the AEHs associated with ae are scheduled, and the **run()** method for each is executed asynchronously

- You do not override run() when defining an AsyncEventHandler subclass, but rather override **handleAsyncEvent()** which is called from run() and does the real work
- run() calls handleAsyncEvent() repeatedly as long as fireCount > 0

Asynchronous event handling example (1)

A hardware device writes sensor data to a pair of double values at locations 0x1000 and 0x1008

The interrupt handler fires a `dataArrival` event when the values are ready



A Sensor object serves as the event handler

- It copies these values into a Position object
- It fires a PositionUpdated event

A Reporter object handles the positionUpdated event

- It displays the values on the console

```
class AsyncEventExample{
    public static void main(String[] args){
        Position p = new Position();
        AsyncEvent dataArrival = new AsyncEvent();
        AsyncEvent positionUpdated = new AsyncEvent();
        Reporter r = new Reporter(p, positionUpdated);
        Sensor s = new Sensor(p, dataArrival, positionUpdated);
    }
}

class Reporter extends AsyncEventHandler{
    final Position pr;

    Reporter(Position p, AsyncEvent positionUpdated){
        pr = p;
        positionUpdated.addHandler(this);
    }

    public void handleAsyncEvent(){
        synchronized(pr){
            System.out.println(pr.x);
            System.out.println(pr.y);
        }
    }
}
```

```
class Sensor extends AsyncEventHandler{
  final RawMemoryFloatAccess rmfa;
  final Position ps;
  final AsyncEvent positionUpdated;

  Sensor(Position p, AsyncEvent dataArrival, AsyncEvent positionUpdated){
    rmfa =
      RawMemoryFloatAccess.createFloatAccess(
        "SRAM", // type of memory
        0x1000, // base address
        16); // size (in bytes)
    ps = p;
    this.positionUpdated = positionUpdated;
    dataArrival.bindTo("InputPortInterrupt");
    dataArrival.addHandler(this);
  }

  public void handleAsyncEvent(){
    synchronized(ps){
      ps.x = rmfa.getDouble(0);
      ps.y = rmfa.getDouble(8);
    }
    positionUpdated.fire();
  }
}
```

Synchronous

- Receiver (“target”) suspends or polls at well-defined point
- Examples: rendezvous, wait / notify, “busy wait”

Asynchronous

- Target is notified “immediately” (no waiting or polling)
- Resumption model
 - Target executes response and then continues from the point of interruption
 - Example: interrupt handler
- Termination model
 - Target executes response but does not continue from the point of interruption



Asynchronous Transfer of Control

- A mechanism whereby a triggering thread (possibly an async event handler) can cause a target thread to branch unconditionally, without any explicit action from the target thread

Uses

- Timing out on a computation
- Awakenng a suspended thread
- Aborting one iteration of a loop
- Aborting a thread

Benefits

- Predictable/bounded response time (avoids latency)

Problems

- Complicated to specify and implement
- Hard to program if control transfers may happen at unpredictable points
 - Triggering thread does not know what state the target thread is in when the ATC is initiated
 - Target thread must be coded carefully in presence of ATC
- May induce overhead even if not used
- *Conflict between desire for ATC to be immediate, and the need for certain code to execute completely*
 - Especially in synchronized code

Semantics

- Triggering thread throws asynchronous exception in target thread
- Control in target thread transfers to handler, or, if none, target thread terminates with unhandled exception

Advantages

- Intuitive conceptual model
- Gives target thread the ability to do cleanup

Potential issues

- Inconsistent state
- Unintended non-termination
- Unintended termination
- Nested ATCs / “competing” exceptions

Semantics

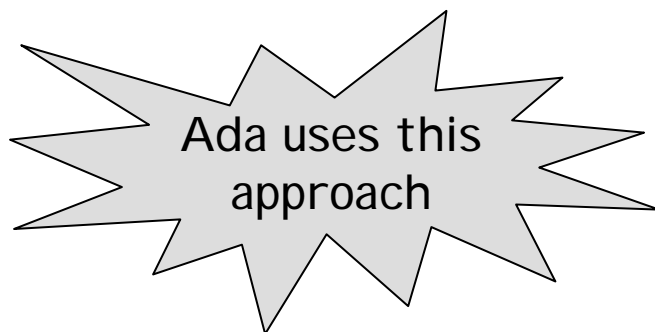
- Thread abort is the basic construct
- Code susceptible to ATC is conceptually an abortable thread
- No exception is thrown

Advantages

- Avoids exception-related anomalies

Potential issues

- Inconsistent state
- How to provide cleanup / prevent “resource leakage”



Methodological principles

- Susceptibility to ATC must be explicit in affected code
- Even if code allows ATC, in some sections ATC must be deferred; in particular, in synchronized code
- Code responding to an ATC does not return to the point where the ATC was initiated

Expressibility principles

- Need mechanism for explicitly triggering an ATC in a target thread
- Need to be able to trigger an ATC based on an external happening, a software event, or a timer
- Need to be able to abort a thread (safely)

Semantic principles

- If ATC modeled through exception handling, need to ensure that the exception is not caught by an unintended handler
- Nested ATCs (for example nested timers) must work properly

Pragmatic principles

- Need idioms for situations such as time-out and thread termination
- Programs that do not use ATC should not incur ATC-related cost

From the viewpoint of the target thread, ATC is modeled by exception handling

- Class `AsynchronouslyInterruptedException` extends `java.lang.InterruptedException`
 - We will use “AIE” as an abbreviation for `AsynchronouslyInterruptedException`
- An ATC is initiated by throwing an instance of AIE, as a result of some action by a triggering thread

Note: allowing a thread to arbitrarily throw an AIE at a target thread would have been problematic

- Target thread may be in code unprepared to catch the exception
 - If exception thrown too early or too late, it propagates and causes target thread to terminate with an unhandled exception
 - A “catch-all” handler may incorrectly handle the exception
 - Semantic mess to define effect of nested AIE propagation
- What if target thread holds a lock when the exception is thrown?
 - Abruptly terminating the synchronized code could lead to inconsistent data (if locks are released) or potential deadlock (if locks are not released)

From the viewpoint of the triggering thread, there are two ways to cause an ATC in a target thread

- Throw a “generic” AIE by invoking `t.interrupt()` where `t` is the target realtime thread
- Throw a specific AIE by invoking `ai e. fire()`
 - The target thread must have indicated its readiness to “accept” `ai e` being thrown asynchronously, else `ai e. fire()` is ignored

ATC only occurs in code that explicitly permits it

- Static indication: method or constructor with a “throws AIE” clause

Some sections of code are *ATC-deferred*

- Synchronized code
- Any method or constructor lacking a “throws AIE” clause

Code that is not ATC-deferred is said to be *asynchronously interruptible* (“AI”)

How a target thread indicates susceptibility to interruption by a specific AIE instance (ai e)

- Express the interruptible code as the run() method of an instance, *i*, of the Interruptible interface
- Invoke ai e.doInterruptible(i)
 - This enables the thread to be interrupted by a triggering thread or AsyncEventHandl er invoking ai e.fire()
 - The AIE object referenced by ai e must be accessible to both threads

Effect of ai e.doInterruptible(i)

- Enable ai e
- Invoke i.run() synchronously
- Disable ai e

Semantics of `ai e. fi re()`

- If `ai e` is enabled
 - If target thread is in an AI section, throw `ai e` immediately
 - Otherwise keep `ai e` as *pending* until execution next reaches an AI section and throw it then
 - “Next” may be either when invoking an AI method or constructor, or when returning (normally or abnormally) to such code
- If `ai e` is disabled, `ai e. fi re()` is a no-op

Semantics of `rt. i nterrupt()` for `Real ti meThread rt`

- Generalizes `Thread. i nterrupt()` to ATC
- A “generic” AIE is made pending on thread `rt`
- This AIE is thrown as soon as `rt` is in an AI section
 - Stays pending even if handled, unless/until a handler calls `cl ear()`

Nested AIE's

- Invocations of `doI nterrupti bl e()` may be nested
- The rules ensure that at most one AIE is pending
 - Precedence is given to the generic AIE and then to the named AIE directed at the shallower scope

```
java.lang.InterruptedException
```

```
public class AsynchronouslyInterruptedException  
    extends java.lang.InterruptedException{  
    public AsynchronouslyInterruptedException();  
    public boolean doInterruptible(Interruptible logic);  
    public boolean disable();  
    public boolean enable();  
    public boolean fire();  
    public boolean clear();  
    public boolean isEnabled();  
    public static AsynchronouslyInterruptedException getGeneric()  
}
```

```
public class Timed extends AIE{  
    public Timed(HighResolutionTime time);  
    public boolean doInterruptible(Interruptible logic);  
    public void resetTime(HighResolutionTime time);  
}
```

```
public interface Interruptible{  
    void run( AsynchronouslyInterruptedException e )  
        throws AsynchronouslyInterruptedException;  
    void interruptAction(AsynchronouslyInterruptedException e );  
}
```

```
new Timed(new RelativeTime( millis, nanos )).doInterruptible(  
  new Interruptible(){  
    public void run(AIE e) throws AIE{  
      ... // code susceptible to asynch interruption  
    }  
  
    public void interruptAction(AIE e){  
      ... // code executed if run interrupted  
    }  
  });
```

Example: use ATC to deliver a function result by successive approximation after a specified time interval

```

abstract class Func{
  abstract double f(double x)
    throws AIE;
  volatile double current;
  // assumes atomic
}

class MyFunc extends Func{
  double f(double x) throws AIE {
    current = ...;
    while(...){
      ... current = ...;
    }
    return current;
  }
}

```

```

class SuccessiveApproximation{
  static boolean finished;
  static double calc(Func func,
                    double arg,
                    long ms){
    double result = 0.0;
    new Timed( new RelativeTime(ms, 0) ).
      doInterruptible(
        new Interruptible(){
          public void run(AIE e) throws AIE{
            result = func.f(arg);
            finished = true;
          }
          public void interruptAction(AIE e){
            result = func.current;
            finished = false;
          }
        }
      ));
    return result;
  }
  public static void main(String[] args){
    MyFunc mf = new MyFunc();
    double answer = calc(mf, 100.0, 1000);
    // run mf.f(100.0) for at most 1 second
    System.out.println(answer);
    System.out.println("calc completed? " +
                      finished );
  }
}

```

Goals

- Terminate a thread “immediately” and unconditionally but safely:
 - Don’t leave shared object in inconsistent state
 - Don’t deadlock

Original Java approach, and its problems

- The deprecated method `t.stop()` throws a `ThreadDeath` exception and unwinds the stack, releasing locks
 - This leaves shared objects inconsistent
 - A handler for `Exception` or `Throwable` will stop the termination
- The recently-deprecated and never-implemented method `t.destroy()` terminates the thread wherever it was, without throwing an exception
 - Thus any locks held by the thread are not released, which will deadlock any other thread that attempts to acquire any of them

Alternative solution in regular Java

- Use `t.interrupt()` and poll in `t.run()` method to see when to return (“cooperative threadicide”)
- Problem: potential latency

RTSJ Approach

- Use `t.interrupt()` and program `t` do its processing in an asynchronously interruptible method

Benefits of RTSJ Approach

- Safe, since the AIE is deferred while the target thread is in synchronized code
- Lower latency than with polling, when interrupted in an asynchronously interrupted section
- Author of target thread decides where the thread can be interrupted for termination

Drawbacks

- Termination not immediate (and indeed, if the thread never reaches an AI method then it will not be terminated)
- Notation can be heavy

General drawback to non-polling approaches

- Programming so that a program is robust in the presence of asynchronous thread termination is difficult

Thread that triggers the termination

```
class Executioner extends Thread{
  Thread victim;
  Executioner(Thread victim){ this.victim=victim; }
  public void run(){
    victim.start();

    ...
    if (...){victim.interrupt();}
    ...
  }
}
```

Style with polling

```
class PollingVictim extends Thread{
  public void run(){
    while (!interrupted()){
      ...
    }
    ... // pre-shutdown actions
  }
}
```

Style with ATC

```
class AbortableVictim extends RealTimeThread{
    private void body() throws AIE{
        ... // abortable here
    }
    public void run(){
        ... // not abortable here
        try{
            body(); // abortable here
        }
        catch( AIE e ){
            ... // response to abort
        }
        finally{
            ... // unconditional pre-shutdown actions
        }
    }
}
```

RTSJ uses lexical rules

- Synchronized code can call an AI method, and the latter can suffer an ATC
 - An AIE is propagated (synchronously) to caller
- ☺ Natural effect of exception-based semantics
- ☺ Reduces latency
- ☹ Complicates ensuring object consistency
- ☹ Issue with inline expansion

Ada uses dynamic rules

- If an operation is abort-deferred, then so is any subprogram that it invokes
- ☺ Preserves object consistency
- ☹ Increases latency

Safety

- ☺ ATC is deferred in synchronized code
- ☺ ATC permission is per-method, via "throws AIE" clause
- ☹ ATC is possible in finally clauses
- ☹ finally clauses in AI code not executed when AIE propagates

Style / expressiveness

- ☺ Low-level mechanisms provide flexibility
- ☹ High-level idioms (e.g. Thread class) help, but basically the style is clumsy

Performance

- ☹ Latency depends on whether called methods are AI
- ☹ Expect distributed overhead, in absence of optimizations / restricted profile

Other

- ☹ Semantics of finally clauses interacts poorly with JSR 166 lock idioms

Semantics

- Ada abort-based model is simpler
- RTSJ has non-standard rules for AIE propagation, other special rules to avoid anomalies

Safety

- Ada has more reliable finalization behavior
 - finally clauses in RTSJ are asynchronously interruptible
 - finally clauses in RTSJ may be skipped in AIE propagation
- But RTSJ has finer-grained (per-method) control over ATC susceptibility, with “ATC deferred” as default
 - RTSJ: safer invocation of legacy code
 - Ada: subprogram’s abort-deferral based on context
- Both ensure no ATC from synchronized / locked code, but at different levels of granularity

Style / expressiveness

- Ada is much more readable, with tailored syntax
- But RTSJ (actually regular Java) allows external unblocking of a suspended thread

Performance

- Latency is function of programming style
- Hard to avoid overhead, even if ATC not used

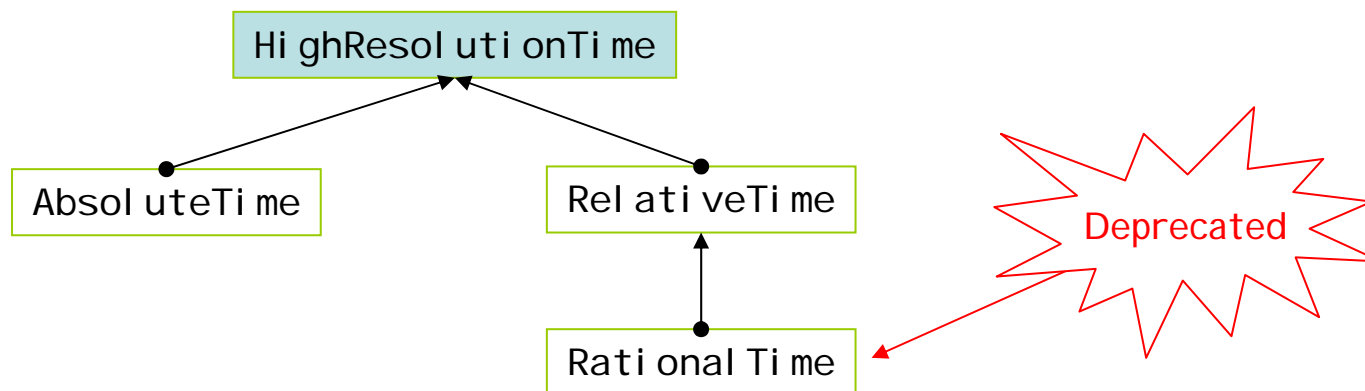
But ATC in general is a “last resort” feature

- Program is more readable, reliable, testable if all control transfers are synchronous
- Use ATC only for specialized circumstances
 - Program termination or “mode change”
 - Contexts where synchronous communication would incur unacceptably high latency

High-resolution time value

- [millis, nanos] where millis is 64 bits (long) and nanos is 32 bits (int)
- Normalized form has both components non-negative or both components non-positive, and nanos in range $-10^6 + 1 .. 10^6 - 1$
- Always based on some Clock (by default the RealtimeClock)

Class hierarchy



Methods for HighResolutionTime

- Conversion to absolute or relative value ("in place" or return new object)
- Comparison (`compareTo`, `equals`)
- "Getters" and "setters"

AbsoluteTime class

- Various constructors, e.g.:
 - `AbsoluteTime(millis, nanos)` constructs value representing point that is [millis, nanos] after the origin of the epoch (00:00:00 GMT, 1/1/1970)
- Methods to associate a new Clock with an AbsoluteTime object
- Arithmetic operations, e.g.:
 - `absTime.add(millis, nanos)` // return new object
 - `absTime.add(millis, nanos, dest)` // update dest
 - `absTime.add(relativeTime)` // add a RelativeTime, return new object
 - `absTime.add(relativeTime, dest)` // update by adding a RelativeTime
- Methods to convert to a relative time based on an explicit clock

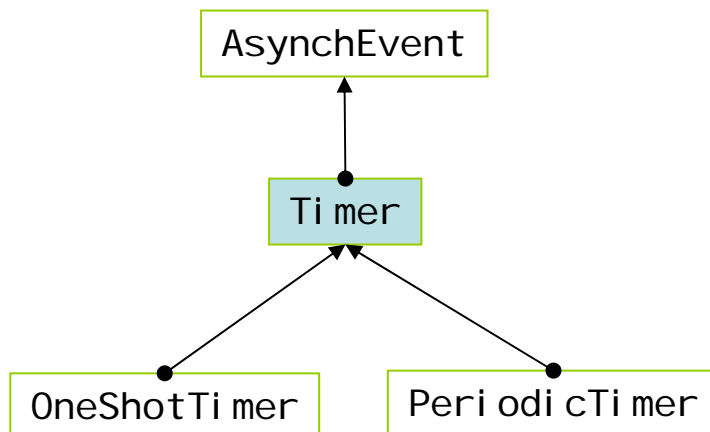
RelativeTime class

- Value represents millis and nanos relative to some AbsoluteTime for some Clock
- “Obvious” set of constructors, e.g. `RelativeTime(millis, nanos)`
- Methods to associate a new Clock with a RelativeTime object
- Arithmetic operations
- Methods to convert to an absolute time based on an explicit clock

The abstract Clock class

- “No-arg” constructor
- Methods
 - `public static Clock getRealTimeClock()`
 - “Singleton” object for a monotonic non-decreasing clock, not necessarily synchronized with external world
 - `public abstract AbsoluteTime getTime()`
 - `public abstract AbsoluteTime getTime(AbsoluteTime time)`
 - Returns the time parameter, into which the current time is copied
 - `public abstract void setResolution(RelativeTime resolution)`

The Timer class hierarchy



Abstract class Timer

- Constructor
 - `protected Timer(HighResolutionTime time, // When to fire event
Clock clock, // Which clock
AsyncEventHandler handler) // Which handler`
- Methods
 - Methods to change timer's state (destroy, disable, enable, reschedule, start, stop)

OneShotTimer class

- Timed AsyncEvent that fires once (subject to enable / disable)
- Constructor is the same as for Timer, plus another one without an explicit Clock parameter

PeriodicTimer class

- AsyncEvent that fires periodically based on its parameters
- Constructors specify interval (RelativeTime) as well as start time, clock (optional), and handler

Example of periodic timer

```
public class Example{
    public void main( String[] args ){
        AsyncEventHandler aeh = new AsyncEventHandler(
            new PriorityParameters( 20 ),
            null, // release parameters
            null, // memory parameters
            null, // memory area
            null, // processing group parameters
            false, // heap
            new Runnable(){
                public void run(){ System.out.println( "Hello" ); }
            }
        )

        PeriodicTimer pt =
            new PeriodicTimer(
                null, // activate at start
                new RelativeTime( 1000, 0 ), // every second
                aeh ); // handler

        pt.start(); // 1st period begins now
        Thread.sleep( 10000 ); // suspend for 10 sec
        pt.destroy();
    }
}
```

Advantages

- Satisfies requirements for nanosecond accessibility, high-resolution
- Supports relative and absolute time
- Accommodates monotonic clock
- Supports different kinds of timers

Disadvantages

- Class hierarchy is suspect
 - `RationalTime` deprecated
 - Occurrences of `ClassCastException` (e.g. for `compareTo` methods in `HighResolutionTime`)
- `Clock` is a weird class
 - Abstract class with a static method returning a singleton instance of one of its subclasses
 - Functionality of real-time clock would be better expressed as a singleton subclass `RealTimeClock` with a static `instance()` method
- Complicated state transition diagram for timers

RawMemoryAccess class

- Can construct a memory area of a given type (e.g. DMA) having a given base and size
- “Peek/poke” of byte, short, int, long (no references)
- Class RawMemoryFloatAccess includes float, double
- Constant BYTE_ORDER in class RealTimeSystem ⇒ endianness
- Example

```
RawMemoryAccess rm =  
    new RawMemoryAccess( "RAM", 0x2000);  
rm.setInt(0x1000, 0xABCDFFFF); // stored in bytes 0x1000 .. 0x1003
```

Physical Memory classes

- Immortal Physical Memory, LTPhysical Memory, VTPhysical Memory
- Use for specialized hardware such as memory-mapped I/O
- Can contain arbitrary objects (not limited to primitive types as in RawMemoryAccess)

Critique

- Captures some of the requirements for low-level programming, but (especially Raw Memory support) goes against the grain of Java

Ada

- Sponsored “top down” effort ⇒ ISO standard + Rationale
- Detailed audit trail (LSNs, AIs, etc.)
- Thorough review (ARG, WG9)
- Highly open process (public briefings, etc.)
- Product evolution based on ISO rules

RTSJ

- Focused “bottom up” volunteer effort ⇒ *de facto* standard
- Java Community Process requires not just the spec but also a Reference Implementation and Technology Compatibility Kit
- Audit trail comprises principally the group’s e-mail messages
- Review was principally internal in RTJEG
- Semi-open process
- Product evolution based on Sun’s JCP rules

Ada

- ☺ Performance (classical stack-based language, queueless lock management)
- ☺ Conservatism (traditional static compile/bind/link)
- ☺ Well-defined semantics (queue placement)
- ☺ Cleaner / simpler approach to ATC
- ☺ Existence of good implementations now
- ☺ Allows but does not require OOP paradigm
- ☹ Market perception

RTSJ

- ☺ Flexibility (multiple schedulers, dynamic loading...)
- ☺ Functionality (admission control, ...)
- ☹ Style may seem complicated to traditional Java programmers
 - Need to pay attention to memory management issues
 - OOP for real-time is not mainstream
- ☹ Performance questions

Specific technical ideas

- Absolute delay (sleep()) method taking an absolute time)
- Scheduling semantics for base scheduler
- Concept of "abort-deferred" regions of code
- Priority ceiling emulation
- Subsets for specialized application areas

Political lessons

- Customers want solutions, not technology
- Beware the culture clash
 - Real-time applications take a static approach to ensure predictability
 - All heap objects are allocated at system startup
 - OOP and garbage collection have not been popular

Challenges

- Sacrificing performance/flexibility for safety (an effect of Garbage Collection) has always been a hard sell to the real-time community

Friends

- “The enemy of my enemy is my friend”
- Cross-fertilization of ideas beneficial to both
 - Many Ada concepts influenced RTSJ and Real-Time Core Extensions
 - Priority Ceiling, ATC, absolute delay, Ravenscar profile
 - RTSJ can serve as model for future Ada work in some areas
 - “On line” feasibility analysis, integrated support for real-time characteristics
- RTSJ-compliant JVM is feasible target for Ada

Foes

- Some organizations are looking at the RTSJ as an alternative to Ada

Peaceful coexistence

- Ada and RTSJ have somewhat different markets
 - Ada: traditional real-time, especially high-integrity
 - RTSJ: organization already committed to Java

Background

- Effort started at same time as the RTSJ, in part due to concerns over licensing issues associated with the Java Community Process
- Derived from Kelvin Nilsen's work on PERC VM and real-time GC

Objectives

- Extend Java platform with real-time capabilities, with services and performance similar to those of commercial RTOSes
- Serve as foundation for higher-level profiles

Functionality similar to RTSJ but from different perspective

- Concurrency model, priority inversion control, allocation contexts, asynchrony, ...
- Underlying decision: clear separation between "Baseline Java" and Core components (e.g., distinct class hierarchies, separate heaps)

Status

- Most recent design document was October 2000
- No implementations
- Main impact is likely to be indirect, via influence on other specs related to real-time Java

Reference Implementation from Timesys

- <http://www.timesys.com>

jRate (Java Real-Time Extension)

- Extends GNU GCJ front-end & run-time library, to support RTSJ
- <http://www.cs.wustl.edu/~corsaro/jRate/>

Ovm (Open VVM)

- Open-Source framework for language run-time systems
- DARPA-sponsored; Purdue, SUNY Oswego, U. of Maryland, DL Tech
- Current thrust: produce JVM compliant with RTSJ
- <http://www.ovmj.org/>

Jamaica VM

- Implementation of RTSJ using F. Siebert's real-time GC algorithm
- <http://www.aicas.com/jamaica.html>
- HIDOORS project: <http://www.hidoors.org/>

FLEX Project (MIT)

- Compiler infrastructure written in Java for Java
- <http://www.flex-compiler.lcs.mit.edu/>

Why consider?

- Sequential Java is secure, with implementation-independent semantics
- Lots of interest in Java technology in systems that have some safety-critical requirements

Some issues

- OOP
- Garbage Collection
- ATC

Approach

- Use RTSJ as basis, subset as required, and perhaps incorporate some ideas from J-Consortium's Core Extensions

Main properties

- No GC
- Limited API
- Compile-time checking where RTSJ has run-time checks

Status

- Initial stages, under The Open Group

RTSJ has addressed the major technical issues concerning Java for real-time applications

- GC unpredictability ⇒ memory areas, NoHeapReal timeThread
- Underspecified concurrency semantics ⇒ base scheduler
- Priority inversions ⇒ monitor control policies
- Issues with stop(), destroy() ⇒ ATC

RTSJ is attracting attention from the real-time research and user communities

Reasons that it may succeed

- Flexibility / generality / extensibility
- “Sexy” technology

Possible obstacles

- Design has “rough edges” (finalized before implementations)
- Performance?
- Style, some local to RTSJ and others due to Java limitations
- Problems with the Java thread model
- User conservatism about OOP

Java thread model

- S. Oaks and H. Wong; *Java Threads (3rd Edition)*; O'Reilly; 2004; ISBN 0-596-00782-5
- D. Lea; *Concurrent Programming in Java (Second Edition)*; Addison-Wesley; 2000; ISBN 0-201-31009-0
- S.J. Hartley, *Concurrent Programming: the Java Programming Language*; Oxford University Press; 1998; ISBN 0-19-511315-2
- A. Holub; *Taming Java Threads*; Apress; 2000; ISBN 1-893115-10-0

Real-time Java: publications

- G. Bollella, J.Gosling, B. Brosgol, P. Dibble, S. Furr, D. Hardin, M. Turnbull, *The Real-Time Specification for Java*, Addison-Wesley; 2000; ISBN 0-201-70323-8
- J-Consortium Specification No. T1-00-01; *Real-Time Core Extensions for the Java Platform*; 2000; Available from www.j-consortium.org
- P. Dibble; *Real-Time Java Platform Programming*; Prentice-Hall; 2002; ISBN 0130282618
- A. Wellings; *Concurrent and Real-Time Programming in Java*; John Wiley & Sons; 2004; ISBN 047084437X

Real-time Java: web sites

- www.rtj.org
 - Real-Time for Java Experts Group
- www.j-consortium.org
 - J-Consortium
- www.nist.gov/rt-java
 - NIST Requirements