# Safe Region-based Memory Management for Java

## Jan Vitek

*with Tian Zhao and Jason Baker*

**PURDUE**
UNIVERSITY

$(\int^3)$

# This talk

- High-level programming languages facilitate software development by abstracting error prone or tedious tasks.

- Memory management is error prone and tedious, but real-time systems sometimes need fine grained control.

- This talk shows how to regain control over memory ... when needed.

- The proposed solution is a hybrid model, a language running mostly on GC but with spurts of manual allocation ... with a new high-level abstraction for Java called *Scopes*.

- Our solution leverages a long line of results from type theory, the key contribution is a statically safe region-based system which is simple and backwards compatible.

- Some related work: MLKit, Cyclone, Ownership types.

# Embedded Real-time Systems

- Real-time embedded systems are central to many applications from avionics, to to automotive industry; they are the largest installed base of microprocessors

- The size of embedded systems is growing steadily; up to multi-million line systems (e.g. DD(X) battleship control)

- Very low level languages (e.g. assembly) are not viable; low-level ones (eg. C ) are barely tolerable; Ada has unfortunately not found the degree of adoption it deserved...

    ... new programming language abstractions for embedded real-time systems are sorely needed.

- Reusability is becoming mandatory, even Boeing changing from the old way (recode from scratch) in favor of COTS
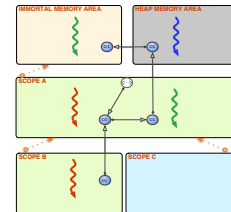
# Automated Memory Management

- Garbage collection algorithms along with memory-safe languages prevent all memory errors.

- The cost is an increase in memory requirements (or a decrease in performance, take you pick) with 2x over the optimal not uncommon.

- But the main issue is pause times. Stop-the-world collectors require the application to pause for >100ms. In a RT system this can mean missing deadlines.

- Recently, real-time collectors have been able to bound pauses to <20ms, this comes at some cost in throughput

- The state of the art in RT GC seems good enough for many RT applications, but not for some hard real-time subsystems

# Region-based allocation

- Region-based allocation is an alternative to GC and manual allocation. It follows from the observation that data exhibits liveness locality (allocation and deallocation times are correlated).

- <u>Idea</u>: create regions which are pools of memory that can be used to store data and bulk deallocated in constant time.

- This simplifies memory management because instead of tracking individual pointers, we deal with entire regions.

- A stack of regions allows to deal with different lifetimes. Parent regions outlive child regions; data is allocated in the region most closely matching its lifetime.

- In itself, region-based allocation is not safe. Programming language techniques must be devised to enforce safety (either run-time or compile-time).

- The <u>Real-time Specification for Java</u> provides region-based MM, but safety is enforced dynamically. This entails runtime overheads and runtime failures.



**PURDUE**
UNIVERSITY

---

# Motivating example: PRISMj



### Mission critical avionics DRE
Boeing, Purdue, UCI, WUSTL
*Route computation, Threat deconfliction algorithms*
*ScanEagle UAV*

| System | K LOCs |
|---|---|
| PRISMJ | 109K |
| FACET EVENT CHANNEL | 15K |
| ZEN CORBA ORB | 179K |
| RTSJ LIBRARIES | 60K |
| CLASSPATH LIBRARIES | 500K |
| OVM VIRTUAL MACHINE | 220K |



| PrismJ avionics controller (app layer) |
|---|
| FACET event channel |
| ZEN  Object Request Broker |
| Real-time Specification for Java (User level implementation) |
| Ovm virtual machine kernel |

kernel boundary

3 rate groups (20, 5, 1Hz)
several hundred RT threads

**Embedded Planet PowerPC 8260**
```
Core at 300 MHz
256 Mb SDRAM
32 Mb FLASH
PC/104 mechanical sized
Embedded Linux
```

**PURDUE**
UNIVERSITY

# SCOPES

# Design Requirements

- The design of scopes abides by the following requirements:

  1. **Static type safety**: *no runtime errors and no runtime checks.*

  2. **GC-safety**: *no interference with the GC so that a thread within a scope may safely preempt the garbage collector and never need block for GC.*

  3. **Bytecode compatibility**: *the output of the compiler should be valid Java bytecode. The only allowed change to the tool chain is the addition of an extra verification pass.*

  4. **Minimal changes to the VM**: *modifications should be limited to the memory subsystem.*

  5. **Source-level backwards compatibility**: *existing Java class should be reusable from scope code.*

# Scope Basics

- Scope objects reify allocation contexts
- Creating a scope allocates a backing store
- Invoking a method on a scope switch allocation context
- Releasing a scope deallocates the backing store

- nb: use of syntactic sugar can be replaced by annotations

```
1   r = new processor;
2   while (b = nextRequest()) {
3       r.getMessage(b,m);
4       m.dispatch();
5       release r;  }
```

# Programming with Scopes

- Scopes have fields, methods, nested classes and scopes.

- Bound classes are always allocated within their enclosing scope.

- Scopes form a tree at runtime, parents always outlive children.

- Invoking a bound class method switches allocation context.

```
1    scope processor {
2        class Unpacker { ...
3            void parse(Bytes b) {...}
4            void write(Message m) {...}
5        }
6        void getMessage(Bytes b, Message m) {
7            Unpacker p = new Unpacker();
8            p.parse(b);
9            p.write(m);
10   }    }
```

# A Scope Pool

- Multiple threads can execute in the same scope and communicate via scope fields
- Releasing a scope deallocates the contents only if no thread is active in it

```
1  scope pool {
2    area[] a;   int cnt;
3    void create(int i) {
4      a=new area[i];
5      for(int j=0;j<i;j++) a[j]=new area;
6    }
7    void run(Message m, Bytes data){
8      int i = cnt++ % a.length;
9      a[i].run(m, data);
10     release a[i];
11   }
12   scope area {
13     void run(Message m, Bytes data){
14       ... // service request
15 } } }
```
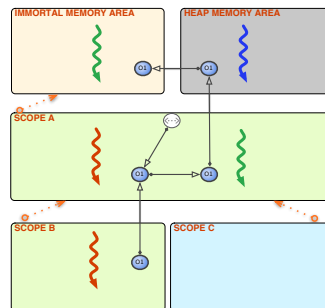
$(\varsigma^3)$

# Implicit Polymorphism

- Classes not lexically enclosed within a scope are **free** classes.
- Free classes are *implicitly* scope-polymorphic.
- Instances of free classes can be allocated within any scope, with the restriction that they not escape their enclosing scope.
- Many library classes can be reused.

```
1  class List {    List next;    Object value; }
2
3  scope user {
4      List l;
5      void create(int i) {
6        List t;
7        for(int j=0;j<i;j++)
8        {  t = new List(); t.next =l; l = t; }
9  }  }
```
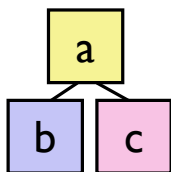
$(\varsigma^3)$

# Permanent Classes

- Permanent classes --annotated with a **top** qualifier-- exist in immortal memory.

- Instances of a permanent class are never deallocated.

- They are allowed to create top-level scopes.

- Permanent classes can be observed from the heap, but cannot refer to heap objects.

- If a thread is spawned in a permanent object, it is guaranteed not to experience GC-interference.
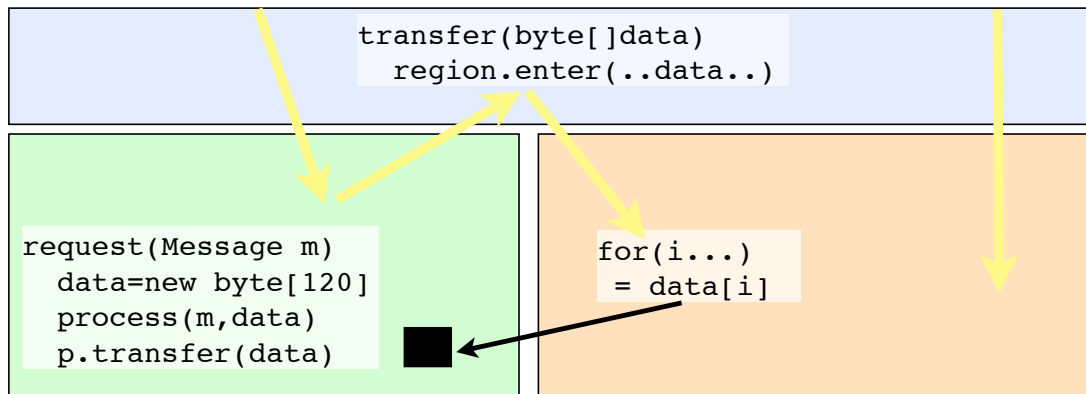
# Crossing Scopes

- Bound classes allocated in parent scopes are visible from subscopes.

- Invoking a method of a bound class changes the allocation context to the parent, when the method returns the allocation context is restored.

- Each thread has a stack of scope, but at any given time only one scope is the effective scope.

```
1   scope a {
2     class Box {
3       c s2;
4       void enter() { s2.mc(); }
5     }
6     scope b { void mb(Box box) { box.enter(); } }
7     scope c { void mc(){...} }
8     void ma() {
9       b s1 = new b;   c s2 = new c;
10      Box box = new Box(); box.s2 = s2;
11      s1.mb(box);
12 }   }
```

# Pipelined computations

```
                    transfer(byte[]data)
                       region.enter(..data..)
```

```
request(Message m)                for(i...)
  data=new byte[120]               = data[i]
  process(m,data)
  p.transfer(data)      ■
```

- LIFO discipline not suited when successive rounds of filtering are needed.

- Pipelined computation can be set up by lending a reference to a sibling region.

- Intuition: a region is pinned while a thread is active within it. If a thread crosses regions, it can safely access its origin scope.

- Safe as long as sibling references are not retained.

# Quasi Linear Types

```
1    scope r {
2      class Bridge {  q qscp;
3          void run(p pscp, q  qscp) {
4                this.qscp=qscp;  pscp.enter(this);
5          }
6          void handoff(borrow byte[] data) {
7                qscp.enter( data);
8      }   }
9      scope p {
10         void enter(Bridge b) {
11            byte[] data=new byte[20];
12            ... // some computation
13            b.handoff( data);
14     }   }
15     scope q {
16         void enter(borrow byte[] data) {
17               ... //use cross-scope reference
18   } }   }
```
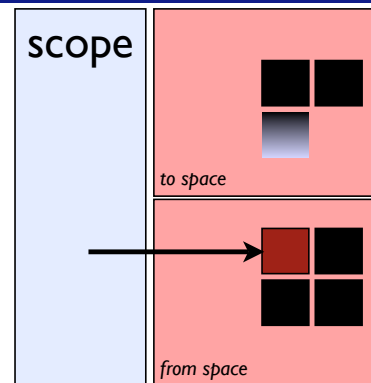
# Quasi Linear Types

- **borrow** -- is a type qualifier that can be used on method arguments and locals.
- The type system ensures a borrowed reference doesn't outlive the method invocation.
- It is safe to access a borrowed reference from any scope.
- Borrowed references cannot be assigned to fields of an object or scope.
- It's safe to modify primitive fields of a borrowed object, but not reference fields.
- Any reference retrieved from a borrowed reference is borrowed as well.

```
                                   class BL{BL next;int i;F f;}
method_m( borrow BL bl ) {
   borrow nbl = bl.next;     // ok
   bl.i++;                   // ok
   bl.f = new F();           // Unsafe! Compile time error.
```
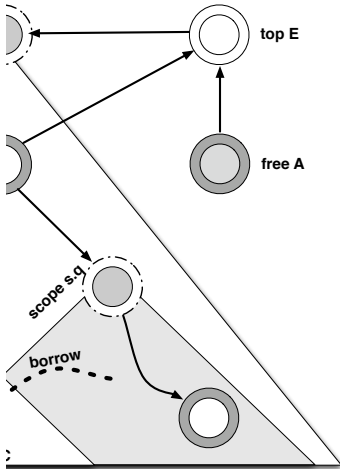
---

# GC-safety

- In a hybrid setting, hard RT code mustn't pause for GC.

- Two scenarios for GC pauses

  - ☑ real-time code allocates in the heap

  - ☑ real-time thread released in midst of GC:
    wait for GC to finish or observe inconsistent data

- Scopes can not refer to the heap thus neither
  can happen

- Do not prevent transitive priority inversion scenarios
    1) a non-rt thread grabs a lock on a permanent object,
    2) GC is triggered,
    3) a RT thread tries to access the locked object.



scope

*to space*

*from space*

- The type system enforces the following invariants:



(1) *Scopes* (`s, s.p, s.q`) may refer to objects within their allocation context and to bound classes of ancestor scopes.

(2) *Bound classes* (e.g. `C, D`) may refer to any co-located object, and to bound classes allocated in a parent scope.

(3) *Free class instances allocated in a scope* (e.g. `B`) may refer to any co-located object.

(4) *Free class instances allocated in the heap* (e.g. `A`) may refer to permanent classes.

(5) *Permanent classes* (e.g. E) may refer to other permanent classes and top-level scopes.

(6) *Borrowed reference* may refer to any bound class.

- Some of the constraints are:
  - ☑ Scopes can not extend other classes
  - ☑ Bound classes can only extend bound classes in the same scope (or free classes)
  - ☑ Free classes can only extend free classes.
  - ☑ Borrowed types can not be stored in fields, or have reference types stored into.
  - ☑ Widening can not cast off the borrowed annotation.
  - ☑ Reference types returned by a borrowed object must be borrowed.
  - ☑ Free types are globally visible
  - ☑ Bound types are visible in the defining scope and subscopes
  - ☑ Scope types are only visible in the enclosing scope
  - ☑ The public interface (field/methods) of scopes and bound types can not contain free classes
  - ☑ Widening is restricted in a number of ways ...

- The scope type system has been formalized as an extension of Igarashi e.a.'s Featherweight Java, and proven sound. The main theorem ensures that after a release, there cannot be any reference into the scope.

$$R ::= \text{scope } S \{ \overline{K\ f}; \overline{M} \}$$
$$L ::= \text{class } C \lhd C \{ \overline{K\ f}; \overline{M} \}$$
$$M ::= K\ m\ (\overline{T\ x}) \{ \text{return } e; \}$$
$$e ::= x \mid v \mid \text{new } K() \mid e.f \mid e.f := e \mid e.m(\overline{e})$$
$$\mid (K)\ e \mid \text{release } e$$
$$v ::= \ell \mid \text{null}$$
$$T ::= K \mid \text{borrow } C$$
$$K ::= C \mid S$$
$$C ::= S.c \mid c$$
$$S ::= S.s \mid \text{top}$$

$$\frac{\Gamma \vdash_{K_t} e : T \quad \textit{fields}(T) = (\overline{K\ f}) \quad T \uparrow K_i = T' \quad (T, K_i) <: (K_t, T') \quad T'\ \textit{viz}\ K_t}{\Gamma \vdash_{K_t} e.f_i : T'} \quad \text{(T-Field)}$$

$$\frac{\Gamma \vdash_{K_t} e : K \quad \textit{fields}(K) = (\overline{K\ f}) \quad \Gamma \vdash_{K_t} e' : K' \quad (K, K_i) <: (K, K_i)}{\Gamma, \Sigma \vdash_{K_t} e.f_i = e' : K'} \quad \text{(T-Upd)}$$

$$\frac{K \neq c \Rightarrow \quad (\textit{scopeof}(K_t) = S \quad K = S.c \ \lor\ K = S.s)}{\Gamma \vdash_{K_t} \text{new } K() : K} \quad \text{(T-New)}$$

$$\frac{\Gamma \vdash_{K_t} e : T_r \quad \Gamma \vdash_{K_t} \overline{e} : \overline{T'} \quad \textit{mtype}(m, T_r) = \overline{T} \to K \quad T_r \uparrow K = T \quad T\ \textit{viz}\ K_t \quad \forall i, (K_t, T'_i) <: (T_r, T_i) \quad (T_r, K) <: (K_t, T)}{\Gamma \vdash_{K_t} e.m(\overline{e}) : T} \quad \text{(T-Invk)}$$

$$\frac{\Gamma \vdash_{K_t} e : K' \quad (K = c \Rightarrow \quad K' = c' \ \lor\ \textit{scopeof}(K') = \textit{scopeof}(K_t))}{\Gamma \vdash_{K_t} (K)\ e : K} \quad \text{(T-Cast)}$$

$$\Gamma \vdash_{K_t} x : \Gamma(x) \quad \text{(T-Var)} \qquad \frac{\Gamma \vdash_{K_t} e : S}{\Gamma \vdash_{K_t} \text{release } e : S} \quad \text{(T-Rel)}$$

**Method typing:**

$$\frac{\overline{x} : \overline{T}, \text{this} : C, \emptyset \vdash_C e : K' \quad K, \overline{T}\ \textit{viz}\ C \quad \textit{override}(m, C, C') \quad (C, K') <: (C, K)}{C \lhd C' \vdash K\ m\ (\overline{T\ x}) \{ \text{return } e; \}} \quad \text{(T-MethC)}$$

$$\frac{\overline{x} : \overline{T}, \text{this} : S, \emptyset \vdash_S e : K' \quad K, \overline{T}\ \textit{viz}\ S \quad (S, K') <: (S, K)}{S \vdash K\ m\ (\overline{T\ x}) \{ \text{return } e; \}} \quad \text{(T-MethS)}$$

**Class typing:**

$$\frac{C \lhd C' \vdash \overline{M} \quad \overline{K}\ \textit{viz}\ C \quad C' = c \ \lor\ (C = S.c \ \land\ C' = S.c')}{\text{class } C \lhd C' \{ \overline{K\ f}; \overline{M} \} \text{ OK}} \quad \text{(T-Class)} \qquad \frac{S \vdash \overline{M} \quad \overline{K}\ \textit{viz}\ S}{\text{scope } S \{ \overline{K\ f}; \overline{M} \} \text{ OK}} \quad \text{(T-Scope)}$$

---

# Conclusions

- High-level languages for real-time systems must support memory-safe, and GC-safe, manual memory management techniques.

- Region-based allocation provides linear time allocation and bulk deallocation.

- Type systems can prevent all memory violations at compile time.

- The scopes abstraction introduced in this talk provides static safety without requiring drastic changes to Java.

- Scopes are simple to use and allow the reuse of existing code.