

## Memory management and programming models

- The choice of memory management affects productivity
- Object-oriented languages naturally hide allocation behind abstraction barriers
  - Taking care of de-allocation manually is more difficult in OO style
- Concurrent algorithms usually emphasize allocation
  - because freshly allocated data is guaranteed to be thread local
  - “transactional” algorithms generate a lot of temporary objects
- ... but garbage collection is a global, costly, operation that introduces unpredictability

## Alternative 1: No Allocation

- If there is no allocation, GC does not run.
  - This approach is used in JavaCard

## Alt 2: Allocation in Scoped Memory

- RTSJ provides scratch pad memory regions which can be used for temporary allocation

‣ Used in deployed systems, but tricky as they can cause exceptions

```
s = new SizeEstimator();
s.reserve(Decrypt.class, 2);
...
shared = new LTMemory(s.getEstimate());
shared.enter(new Run(){ public void run(){
    ...d1 = new Decrypt() ...
}});
```

## Alt 3: Real-time Garbage Collection

- There are three main families of RTGC implementations

- **Work-based**

‣ *Aicas JamaicaVM*

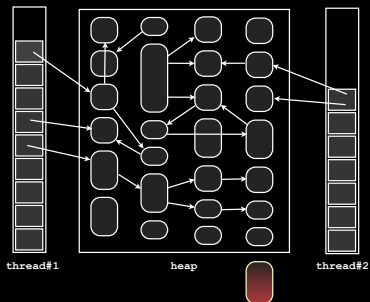
- **Time-triggered, periodic**

‣ *IBM Websphere*

- **Time-triggered, slack**

‣ *SUN Java RealTime System*

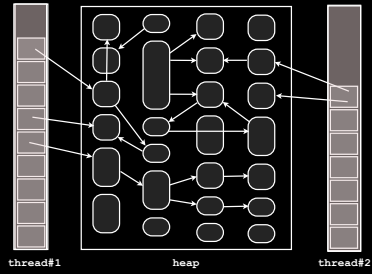
## Garbage Collection



### Phases

- Mutation
- Stop-the-world
- Root scanning
- Marking
- Sweeping
- Compaction

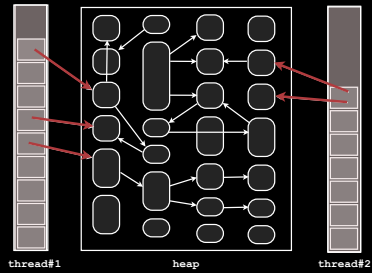
## Garbage Collection



### Phases

- Mutation
- Stop-the-world
- Root scanning
- Marking
- Sweeping
- Compaction

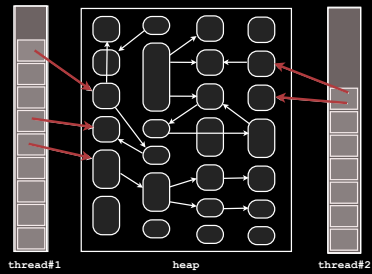
## Garbage Collection



### Phases

- Mutation
- Stop-the-world
- Root scanning
- Marking
- Sweeping
- Compaction

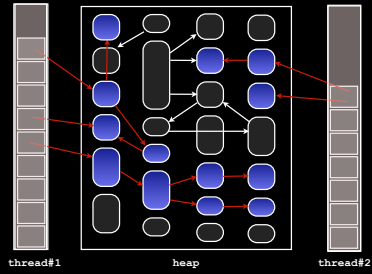
## Garbage Collection



### Phases

- Mutation
- Stop-the-world
- Root scanning
- Marking
- Sweeping
- Compaction

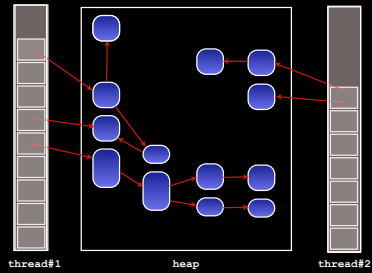
## Garbage Collection



### Phases

- Mutation
- Stop-the-world
- Root scanning
- Marking
- Sweeping
- Compaction

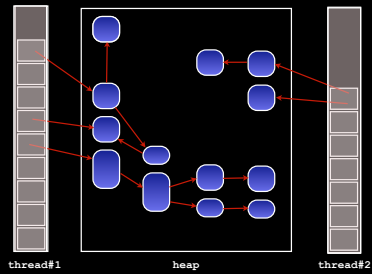
## Garbage Collection



### Phases

- Mutation
- Stop-the-world
- Root scanning
- Marking
- Sweeping
- Compaction

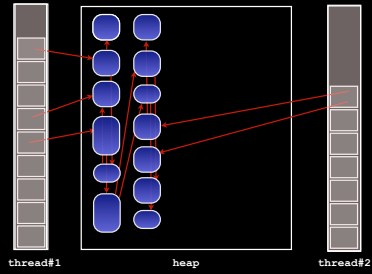
## Garbage Collection



### Phases

- Mutation
- Stop-the-world
- Root scanning
- Marking
- Sweeping
- Compaction

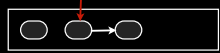
## Garbage Collection



### Phases

- Mutation
- Stop-the-world
- Root scanning
- Marking
- Sweeping
- Compaction

## Incrementalizing marking



Collector marks object



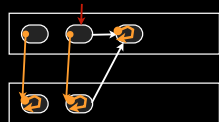
Application updates  
reference field



Compiler inserted  
write barrier marks object

## Incrementalizing compaction

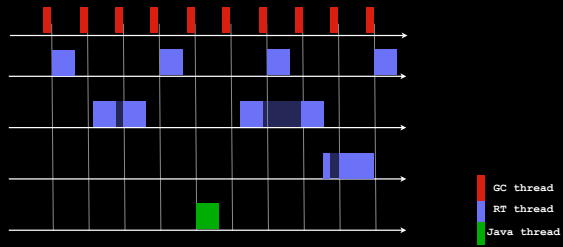
- Forwarding pointers refer to the current version of objects
- Every access must start with a dereference



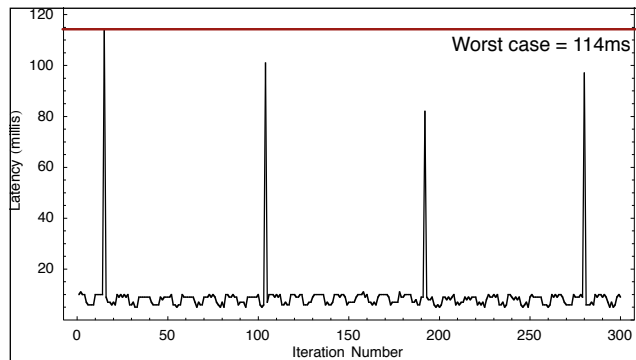
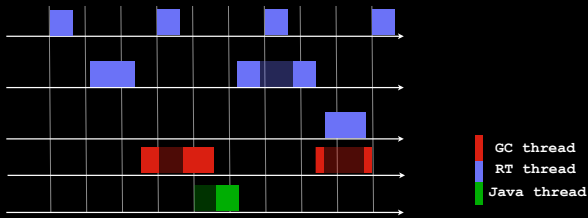
original

copy

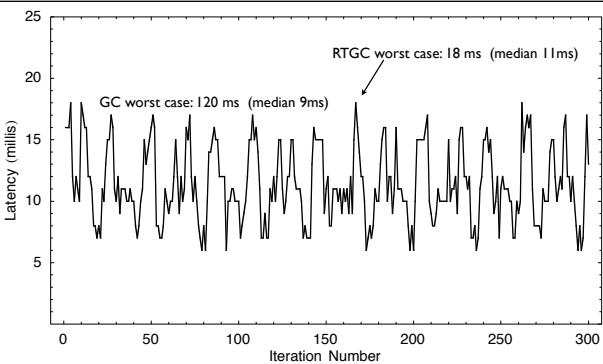
## Time-based GC Scheduling



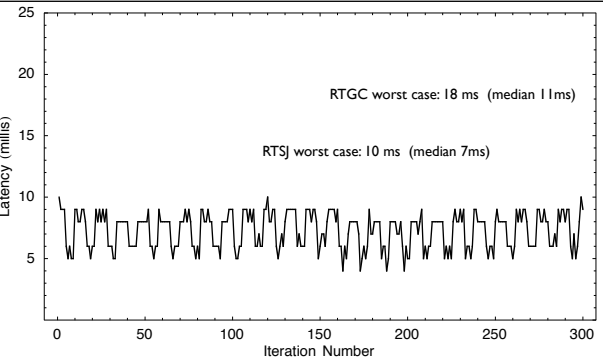
## Slack-based GC Scheduling



- GC pauses cause the collision detector to miss deadlines...
- and this is not a particularly hard problem should support KHz periods



CD with periodic RTGC



Slack-based GC

# Scheduling GC

- Basic parameters needed to schedule a set of task that rely on a real-time garbage collection

$C_i$	[seconds]	computation time	task
$T_i$	[seconds]	period	task
$A_i$	[bytes]	allocation	task
$G_i$	[seconds]	GC work generated	task
$H$	[bytes]	heap size	system
$L_{\max}$	[bytes]	live memory	system
$T_{gc}$	[seconds]	GC cycle duration (period)	system
$G_0$	[seconds]	GC cycle overhead	system

## Traditional response time analysis

32

- C and T are the traditional task WCET and deadline
- Computing the response time analysis can be done as usual by:

$$R_i = C_i + \sum_{j=1}^{i-1} \left( \left\lceil \frac{R_i}{T_j} \right\rceil C_j \right)$$

## Schedulability analysis with GC

33

- The period of the GC task  $T_{gc}$  must be chosen so that it larger than a GC cycle (from start collection to finish)
- The constant  $G_i$  is the upper bound on the amount of GC work that a user task  $i$  can generate in release
- The constant  $G_0$  is the task independent work performed during each GC cycle
- The maximum heap size  $L$  is a chosen by the user
- The maximum amount of live memory  $L_{max}$  depend on the program
- The upper bound on allocation  $A_i$  per release of task  $i$

## Schedulability tests

34

- In the presence of GC two new tests are added to the schedulability equations:
- T1 *The mutator tasks meet their deadline (modified)*
- T2 *The system does not run out of memory*
- T3 *The GC task meets its deadline and keeps with up mutator tasks*



## T2 Memory test

35

- If an object is freed during one GC cycle, the earliest time it can be found and freed is during the following GC cycle
- Assume that we compute  $A_{max}$ , the maximum amount allocated by all mutators during one GC cycle
- The maximum of floating garbage is  $A_{max}$ , plus the maximum allocated in the cycle  $A_{max}$  and the maximum live memory  $L_{max}$  must be less than the heap size  $H$

$$A_{max} \leq \frac{H - L_{max}}{2}$$

## T2 Memory test

36

- If we assume that the GC period is a multiple of the hyperperiod of the mutator tasks, then we can compute:

$$G_{max} = G_0 + \sum_{i=1}^n \frac{T_{gc}}{T_i} G_i$$

$$A_{max} = \sum_{i=1}^n \frac{T_{gc}}{T_i} A_i$$

## T3 GC test for Slack-scheduled

37

- For a slack-scheduled RTGC, the response time of the GC is the maximal amount of GC work that has to be performed per cycle plus the sum of the interruptions of all other tasks.
- The response time of the mutator tasks (T1) is computed as usual

$$R_{gc} = G_{max} + \sum_{i=1}^n \left( \left\lceil \frac{R_{gc}}{T_i} \right\rceil \cdot C_i \right)$$

- A periodic GC is run according to a pattern, e.g.  
CMMCM...
  - For any window of time  $t$  the minimum mutator utilization is the least ratio of time available to the mutator threads, written  $mmu(t)$
  - $mcu(t) = 1 - mmu(t)$  is the ratio used by the collector during that window

$$R_i = C_i + \sum_{j=1}^{i-1} \left( \left\lceil \frac{R_i}{T_j} \right\rceil \cdot C_j \right) + (1 - mmu(R_i)) \cdot R_i$$

- The response time for the GC is computed by finding:

Let  $R_{gc}$  be the smallest  $t$  such that  $t \cdot mcu(t) \geq G_{max}$  and  $t \leq T_{gc}$

