

# A RISC-V JAVA UPDATE

## RUNNING FULL JAVA APPLICATIONS ON FPGA-BASED RISC-V CORES WITH JIKESRVM

**Martin Maas** Krste Asanovic John Kubiatowicz

7<sup>th</sup> RISC-V Workshop, November 28, 2017



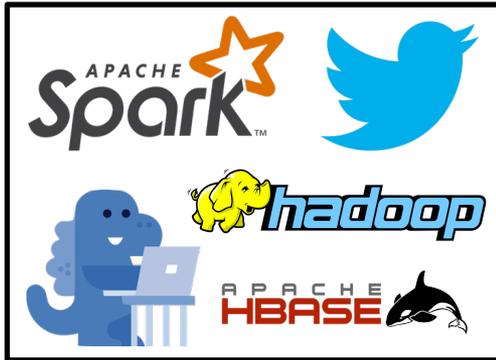
**Berkeley  
Architecture  
Research**



# Managed Languages

# Managed Languages

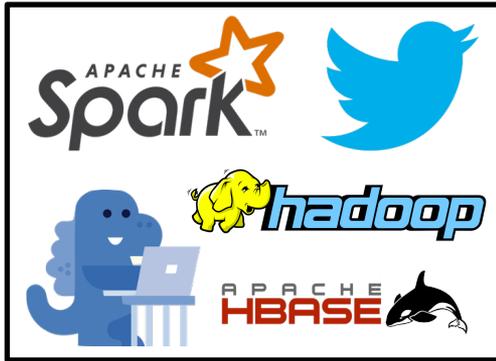
Java, PHP, C#,  
Python, Scala



Servers

# Managed Languages

Java, PHP, C#,  
Python, Scala



JavaScript,  
WebAssembly

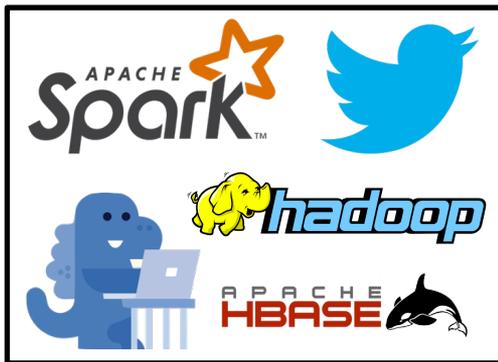


Servers

Web Browser

# Managed Languages

Java, PHP, C#,  
Python, Scala



Servers

JavaScript,  
WebAssembly



Web Browser

Java, Swift,  
Objective-C



Mobile

# Java on RISC-V

# Java on RISC-V

**OpenJDK/Hotspot JVM**



**High-performance production JVM**

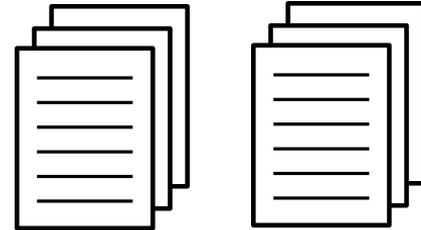
# Java on RISC-V

## OpenJDK/Hotspot JVM



High-performance production JVM

## Jikes Research VM



Easy-to-modify research JVM

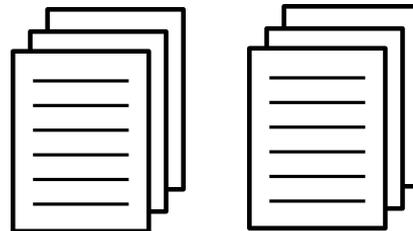
# Java on RISC-V

## OpenJDK/Hotspot JVM



High-performance production JVM

## Jikes Research VM



Easy-to-modify research JVM

# Talk Outline

# Talk Outline

## 1. Running JikesRVM on Rocket Chip

Executing JikesRVM on FPGA-based RISC-V hardware

# Talk Outline

## 1. Running JikesRVM on Rocket Chip

Executing JikesRVM on FPGA-based RISC-V hardware

## 2. Managed-Language Use Cases

New research that is enabled by this infrastructure

# Talk Outline

## 1. Running JikesRVM on Rocket Chip

Executing JikesRVM on FPGA-based RISC-V hardware

## 2. Managed-Language Use Cases

New research that is enabled by this infrastructure

## 3. The State of Java on RISC-V

Progress, Challenges and Announcements

# PART I

- 1. Running JikesRVM on Rocket Chip**  
Executing JikesRVM on FPGA-based RISC-V hardware
- 2. Managed-Language Use Cases**  
New research that is enabled by this infrastructure
- 3. The State of Java on RISC-V**  
Progress, Challenges and Announcements

# JikesRVM on RISC-V

- Runs **full JDK6 applications**, including the Dacapo benchmark suite (no JDK7)
- Passes **JikesRVM core test suite**
- **15,000 lines of code in 86 files** to port the non-optimizing **baseline compiler**

# Porting The Jikes Research VM

## CARRV-2017 Workshop Paper

### Full-System Simulation of Java Workloads with RISC-V and the Jikes Research Virtual Machine

Martin Maas  
University of California, Berkeley  
maas@eecs.berkeley.edu

Krste Asanović  
University of California, Berkeley  
krste@eecs.berkeley.edu

John Kubiawicz  
University of California, Berkeley  
kubitron@eecs.berkeley.edu

```
0.652000] disk [generic-blkdev] of loaded; 381808 sectors, 1 tags, 16 max request length
[ 0.652000] HWGC: Initializing driver, add device with 'mknod /dev/hwgc0 c 254 0'
[ 0.656000] HWGC: IO Region ffffffff44007000
[ 0.676000] VFS: Mounted root (ext2 filesystem) readonly on device 253:0.
[ 0.684000] devtmpfs: mounted
[ 0.688000] Freeing unused kernel memory: 88K
[ 0.692000] This architecture does not have kernel memory protection.
INIT: version 2.88 booting
[ 1.360000] random: fast init done
[ 3.144000] EXT2-fs (generic-blkdev): warning: mounting unchecked fs, running e2fsck is recommended
bootlogd: cannot find console device 4:0 under /dev
hwclock: can't open '/dev/misc/rtc': No such file or directory
Thu Sep 21 17:11:37 UTC 2017
hwclock: can't open '/dev/misc/rtc': No such file or directory
INIT: Entering runlevel: 5
=====
-- JIKES RVM TEST ENVIRONMENT --
=====
chmod: driver-test: No such file or directory
=====
bash: cannot set terminal process group (138): Inappropriate ioctl for device
bash: no job control in this shell
bash-4.4#
bash-4.4#
bash-4.4# ./rvm -X:verboseBoot=10 CARRV
[ 18.260000] random: crng init done
```

```
[ 18.260000] random: crng init done
Booting
Setting up current RVMThread
Doing thread initialization
Setting up memory manager: bootrecord = 0x0000000031000018
Initializing baseline compiler options to defaults
Fetching command-line arguments
Early stage processing of command line
Collector processing rest of boot options
Initializing bootstrap class loader: ./jksvm.jar:./rvmrt.jar
Running various class initializers
running class initializer for java.util.WeakHashMap
invoking method < BootstrapCL, Ljava/util/WeakHashMap; >.<clinit> ()V
running class initializer for org.jikesrvm.classloader.Atom$InternedStrings
invoking method < BootstrapCL, Lorg/jikesrvm/classloader/Atom$InternedStrings; >.<clinit> ()V
running class initializer for gnu.classpath.SystemProperties
invoking method < BootstrapCL, Lgnu/classpath/SystemProperties; >.<clinit> ()V
running class initializer for java.lang.Throwable$StaticData
invoking method < BootstrapCL, Ljava/lang/Throwable$StaticData; >.<clinit> ()V
running class initializer for java.lang.Runtime
invoking method < BootstrapCL, Ljava/lang/Runtime; >.<clinit> ()V
running class initializer for java.lang.System
invoking method < BootstrapCL, Ljava/lang/System; >.<clinit> ()V
running class initializer for sun.misc.Unsafe
invoking method < BootstrapCL, Lsun/misc/Unsafe; >.<clinit> ()V
running class initializer for java.lang.Character
```



# FPGA Performance Results

Benchmarks	Instructions (B)	Simulated Time (s)
avro	118.0	311.8
luindex	47.4	103.5
lusearch	263.5	597.2
pmd	158.5	346.8
sunflow	504.8	1,352.9
xalan	190.8	466.4

**Default input sizes, >1 trillion instructions**

# PART II

1. **Running JikesRVM on Rocket Chip**  
Executing JikesRVM on FPGA-based RISC-V hardware
2. **Managed-Language Use Cases**  
New research that is enabled by this infrastructure
3. **The State of Java on RISC-V**  
Progress, Challenges and Announcements

# Wake Up and Smell the Coffee: Evaluation Methodology for the 21st Century

CACM  
08/2008

By Stephen M. Blackburn, Kathryn S. McKinley, Robin Garner, Chris Hoffmann, Asjad M. Khan, Rotem Bentzur, Amer Diwan, Daniel Feinberg, Daniel Frampton, Samuel Z. Guyer, Martin Hirzel, Antony Hosking, Maria Jump, Han Lee, J. Eliot B. Moss, Aashish Phansalkar, Darko Stefanović, Thomas VanDrunen, Daniel von Dincklage, and Ben Wiedermann

## Abstract

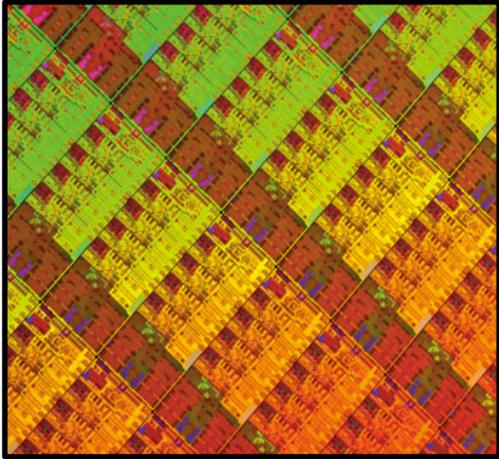
Evaluation methodology underpins all innovation in experimental computer science. It requires relevant *workloads*, appropriate *experimental design*, and *rigorous analysis*.

Unfortunately, methodology is not keeping pace with the changes in our field. The rise of managed languages such as Java, C#, and Ruby in the past decade and the imminent rise of commodity multicore architectures for the next decade pose new methodological challenges that are not yet widely understood. This paper explores the consequences of our collective inattention to methodology on innovation,

Many developers today choose managed languages, which provide: (1) memory and type safety, (2) automatic memory management, (3) dynamic code execution, and (4) well-defined boundaries between type-safe and unsafe code (e.g., JNI and Pinvoke). Many such languages are also object-oriented. Managed languages include Java, C#, Python, and Ruby. C and C++ are not managed languages; they are compiled-ahead-of-time, not garbage collected, and unsafe. Unfortunately, managed languages add at least three new degrees of freedom to experimental evaluation: (1) a *space-time trade-off* due to garbage collection, in which heap size is a control vari-

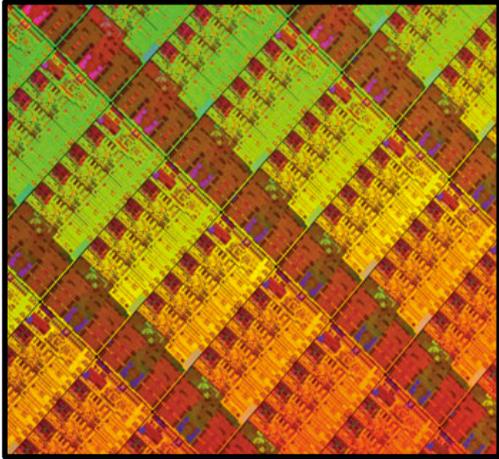
# Managed Language Challenges

# Managed Language Challenges



Long-Running  
on Many Cores

# Managed Language Challenges

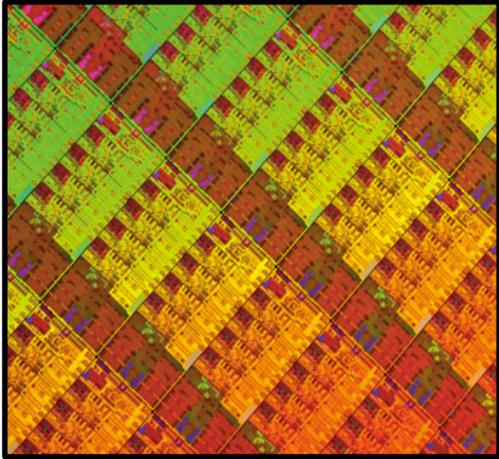


Long-Running  
on Many Cores



Concurrent  
Tasks (GC, JIT)

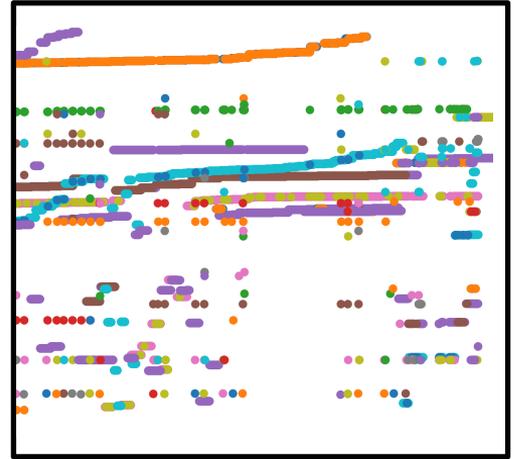
# Managed Language Challenges



Long-Running  
on Many Cores



Concurrent  
Tasks (GC, JIT)



Fine-grained  
Interactions

# Limitations of Simulators

# Limitations of Simulators



Qemu

## High-performance Emulation

Cannot account for fine-grained details  
(e.g., barrier delays of ~10 cycles)

# Limitations of Simulators



Qemu

## High-performance Emulation

Cannot account for fine-grained details  
(e.g., barrier delays of ~10 cycles)

---



## Cycle-accurate Simulation

Too slow to run large-scale Java workloads

# Limitations of Simulators



Qemu

## High-performance Emulation

Cannot account for fine-grained details  
(e.g., barrier delays of ~10 cycles)



## Cycle-accurate Simulation

Too slow to run large-scale Java workloads

Realism

# Limitations of Simulators

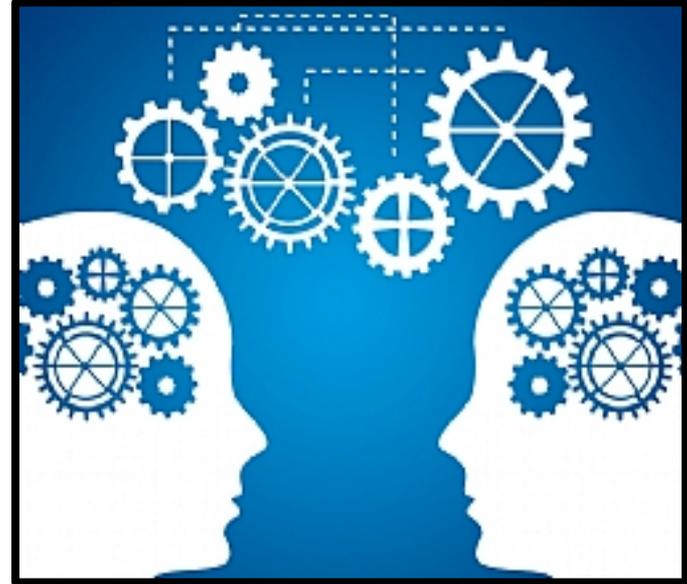


Realism

# Limitations of Simulators



Realism



Industry Adoption



**RISC-V**



**Run managed workloads on real RISC-V hardware in FPGA-based simulation to enable modifying the entire stack**

# Hardware-Software Co-Design

## Grail Quest: A New Proposal for Hardware-assisted Garbage Collection

Martin Maas, Krste Asanović, John Kubiatowicz  
University of California, Berkeley

### ABSTRACT

Many big data systems are written in garbage-collected languages and GC has a substantial impact on throughput, responsiveness and predictability of these systems. However, despite decades of research, there is still no “Holy Grail” of GC: a collector with no noticeable impact, even on real-time applications. Such a collector needs to achieve freedom from pauses, high GC throughput and good memory utilization, without slowing down application threads or using substantial amounts of compute resources.

In this paper, we propose a step towards this elusive goal by revisiting the old idea of moving GC into hardware. We discuss the trends that make it the perfect time to revisit this approach and present the design of a hardware-assisted GC that aims to reconcile the conflicting goals. Our system is work in progress and we discuss design choices, trade-offs and open questions.

### 1. INTRODUCTION

A substantial portion of big data frameworks – and large-scale distributed workloads in general – are written in languages with Garbage Collection (GC), such as Java, Scala, Python or R. Due to its importance for a wide range of workloads, Garbage Collection has seen tremendous research efforts for over 50 years. Yet, we arguably still don’t have what has been called the “Holy Grail” of GC [1]: a pause-free collector that achieves high memory utilization and high GC throughput (i.e., sustaining high allocation rates), preferably without a large resource cost for the application.

Many recent GC innovations have focused on the first three goals, and modern GCs can be made effectively pause-free at the cost of slowing down application threads and using a substantial amount of resources. Moreover, these approaches oftentimes ignore another factor that is very important in warehouse-scale computers: energy consumption. Previous work [2] has shown that GC can account for up to 25% of energy and 40% of execution time in common workloads (10% on average). Worse, as big data systems are processing ever larger heaps, these numbers will likely increase.

We believe that we can reconcile low pause times and energy efficiency by revisiting the old idea of moving GC into hardware. Our goal is to build a GC that simultaneously achieves high GC throughput, good memory utilization, pause times indistinguishable from LLC misses and energy efficiency. We build on an algorithm that performs well on the first three criteria but is resource-intensive [3]. Our key insight is that this algorithm can be made energy efficient by moving it into hardware, combined with some algorithmic changes.

We are not the first to propose hardware support for GC [3–7]. However, none of these schemes has been widely adopted. We believe that there are three reasons

Garbage-collected languages are widely used, but they are rarely the only workload on a system. Systems designed for specific languages mostly lost out to general-purpose cores, partly due to Moore’s law and economies of scale allowing these cores to quickly outperform the specialized ones. This is changing today, as the slow-down of Moore’s law makes it more attractive to use chip area for accelerators to improve common workloads, such as garbage-collected applications.

Most garbage-collected workloads run on servers (note that there are exceptions, such as Android applications). Servers traditionally use commodity CPUs and the bar for adding hardware-support into such a chip is very high (take Hardware Transactional Memory as an example). However, this is changing: cloud hardware and rack-scale machines in general are expected to switch to custom SoCs, which could easily incorporate IP to improve GC performance and efficiency.

Many proposals were very invasive and would require re-architecting of the memory system or other components [5, 7, 8]. We believe an approach has to be relatively non-invasive to be adopted. The current trend to accelerators and processing near memory may make it easier to adopt similar techniques for GC without substantial modifications to the architecture.

We therefore think it is time to revisit hardware-assisted GC. In contrast to many previous schemes, we focus on making our design sufficiently non-invasive to incorporate into a server or mobile SoC. This requires isolating the GC logic into a small number of IP blocks and limiting changes outside these blocks to a minimum.

In this paper, we describe our proposed design. It exploits two insights: First, overlapses of concurrent GCs stem from a large number of small but frequent slow-downs spread throughout the execution of the program. We move the culprits (primarily barriers) into hardware to alleviate their impact and allow out-of-order cores to speculate over them. Second, the most resource-intensive phases of a GC (marking and relocation) are a poor fit for general-purpose cores. We move them into accelerators close to DRAM, to save power and area.

### 2. BACKGROUND

An extensive body of work has been published on GC. Jones and Lins [9] provide a general introduction.

There are two fundamental GC strategies: tracing and reference counting. Tracing collectors start from a set of roots (such as static or stack variables), perform a

## Grail Quest: A New Proposal for Hardware-Assisted Garbage Collection

## 6th Workshop on Architectures and Systems for Big Data (ASBD '16), Seoul, Korea, June 2016

# Motivating Example

We found that the distortion introduced [by the method] unacceptably large and erratic. For example, with the GenMS collector, the [benchmark] reports a 12% to 33% increase in runtime versus running [without].

Modifiable hardware enables fine-grained measurement and injection of language-level data **without** disturbing the application performance

## Quantifying the Performance of Garbage Collection vs. Explicit Memory Management

Matthew Hertz<sup>\*</sup>  
Cornell University  
Ithaca, NY 14850  
matt@hertz@cornell.edu

Emery D. Berger  
Dept. of Computer Science  
University of Massachusetts Amherst  
Amherst, MA 01002  
emery@cs.umass.edu

### ABSTRACT

Garbage collection yields numerous software engineering benefits, but its quantitative impact on performance remains elusive. One can compare the cost of conservative garbage collection to explicit memory management in C/C++ programs by looking at an appropriate collector. This kind of direct comparison is not possible for languages designed for garbage collection (e.g., Java), because programs in these languages naturally do not contain calls to `free`. Thus, the actual gap between the time and space performance of explicit memory management and precise copying garbage collection remains unknown.

We introduce a novel experimental methodology that lets us quantify the performance of precise garbage collection versus explicit memory management. Our system allows us to test multithreaded Java programs as if they used explicit memory management by relying on oracles to insert calls to `free`. These oracles are generated from profile information gathered in earlier application runs. By executing inside an architecturally-extended simulator, the “oracle” memory manager eliminates the effects of coexisting an oracle while measuring the costs of calling `malloc` and `free`. We evaluate two different oracles: a trace-based oracle that aggressively frees objects immediately after their last use, and a reachability-based oracle that conservatively frees objects just after they are last reachable. These oracles open the range of possible placement of explicit deallocation calls.

We compare explicit memory management to both copying and non-copying garbage collectors across a range of benchmarks using the explicit memory manager, and present all non-simulated runs that lend further validity to our results. These results quantify the time-space tradeoff of garbage collection: with four times as much memory, an Apple-style generational collector with a non-copying mutator matches the performance of reachability-based explicit memory management. With only three times as much memory, the collector runs at average 7% slower than explicit memory management. However, with only twice as much memory, garbage collection degrades performance by nearly 70%. When

<sup>\*</sup>Work performed at the University of Massachusetts Amherst.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted, on condition that the copies are not made for commercial sale and that the fee code on this page for copying beyond that permitted by Article 17, Section 107 of the Copyright Act of 1976 is paid to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923. Copyright 2009 ACM 978-1-59593-010-0/09/05...\$5.00

physical memory is scarce, paging causes garbage collection to run an order of magnitude slower than explicit memory management.

### Categories and Subject Descriptors

D.1.1 [Programming Languages]: Dynamic storage management; D.1.4 [Processors]: Memory management (garbage collection)

### General Terms

Experimentation, Measurement, Performance

### Keywords

oracle memory management, garbage collection, explicit memory management, performance analysis, time-space tradeoff, throughput, paging

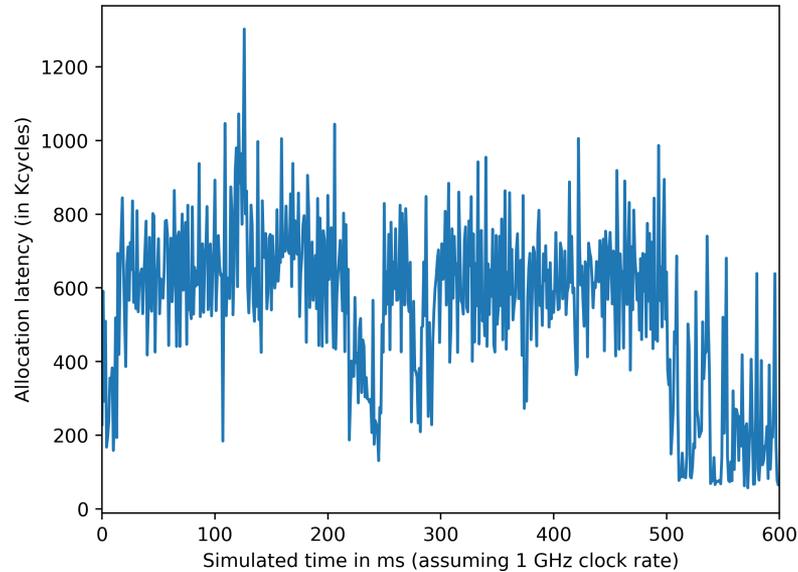
### 1. Introduction

Garbage collection, or automatic memory management, provides significant software engineering benefits over explicit memory management. For example, garbage collection frees programmers from the burden of memory management, eliminates most memory leaks, and improves modularity, while preventing accidental memory overwrites (“dangling pointers”) [55, 59]. Because of these advantages, garbage collection has been represented as a fixture of a number of mainstream programming languages.

Garbage collection can improve programmer productivity [48], but its impact on performance is difficult to quantify. Previous researchers have measured the runtime performance and space impact of conservative, non-copying garbage collection in C and C++ programs [19, 62], for these programs, comparing the performance of explicit memory management to conservative garbage collection is a matter of looking in a library like the Boehm-Demers-Weiser collector [14]. Unfortunately, measuring the performance trade-off in languages designed for garbage collection is not so straightforward. Because programs written in these languages do not explicitly deallocate objects, one cannot simply replace garbage collection with an explicit memory manager. Extrapolating the results of studies with conservative collectors is impossible because precise, deallocating garbage collectors (available only for garbage-collected languages) consistently report more efficient, time-space trading garbage collection [10, 12].

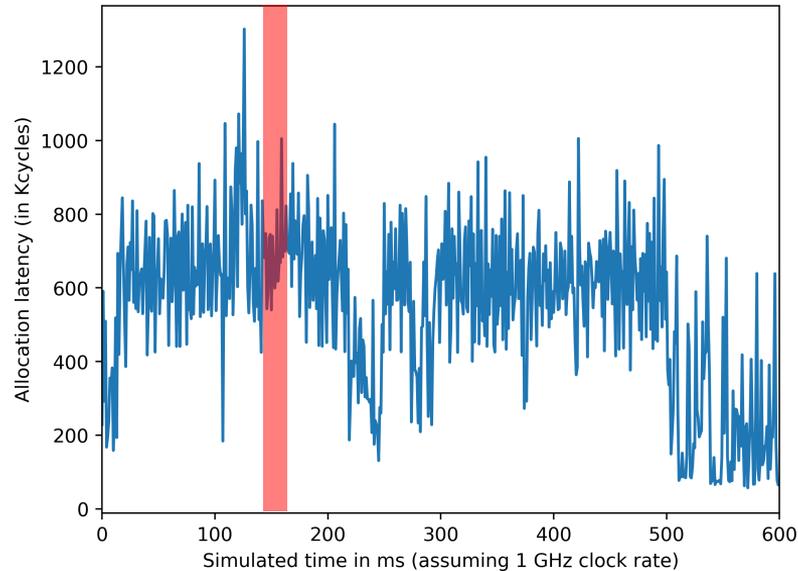
It is possible to measure the costs of garbage collection activity (e.g., tracing and copying) [10, 20, 30, 36, 56] but it is impossible to measure garbage collection activity without measuring the performance of garbage collection. Garbage collection affects application behavior both by visiting and rearranging memory. It also degrades locality, especially when physical memory is scarce [81]. Substituting the costs of garbage collection also ignores the improved locality that explicit memory managers can provide by immediately recycling just-freed memory [55, 58, 57, 58]. For all these reasons, the costs of precise,

# Memory Allocation Latency



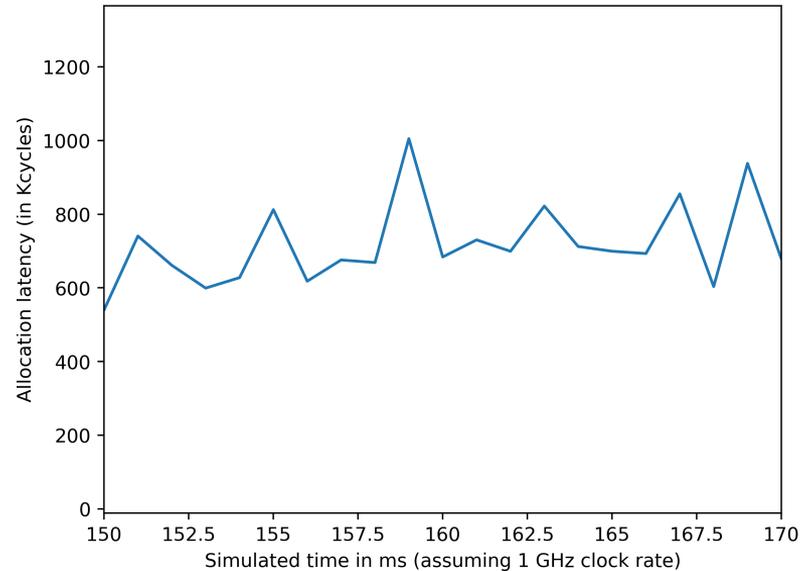
**Sampling Rate: 1 KHz**

# Memory Allocation Latency



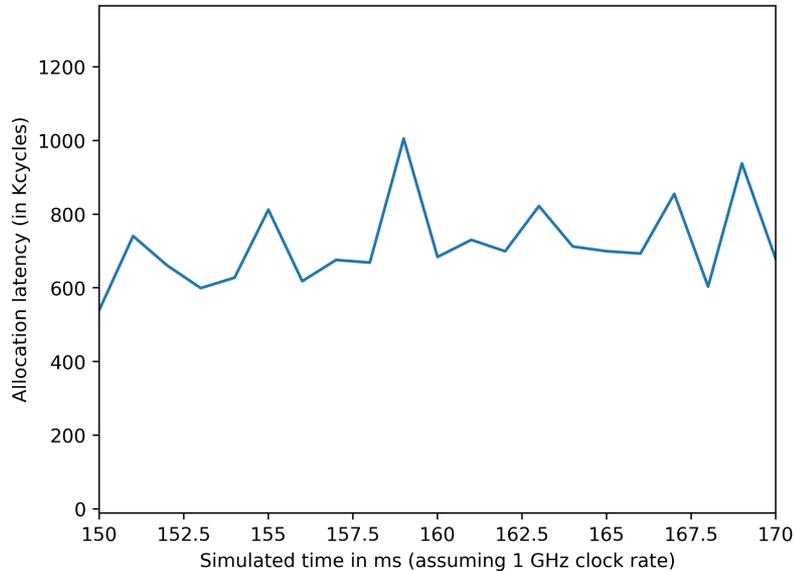
**Sampling Rate: 1 KHz**

# Memory Allocation Latency

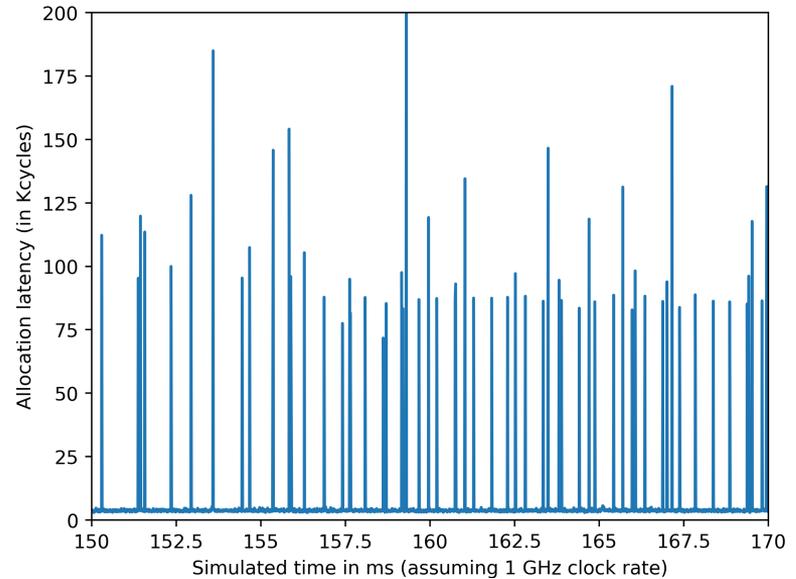


**Sampling Rate: 1 KHz**

# Memory Allocation Latency

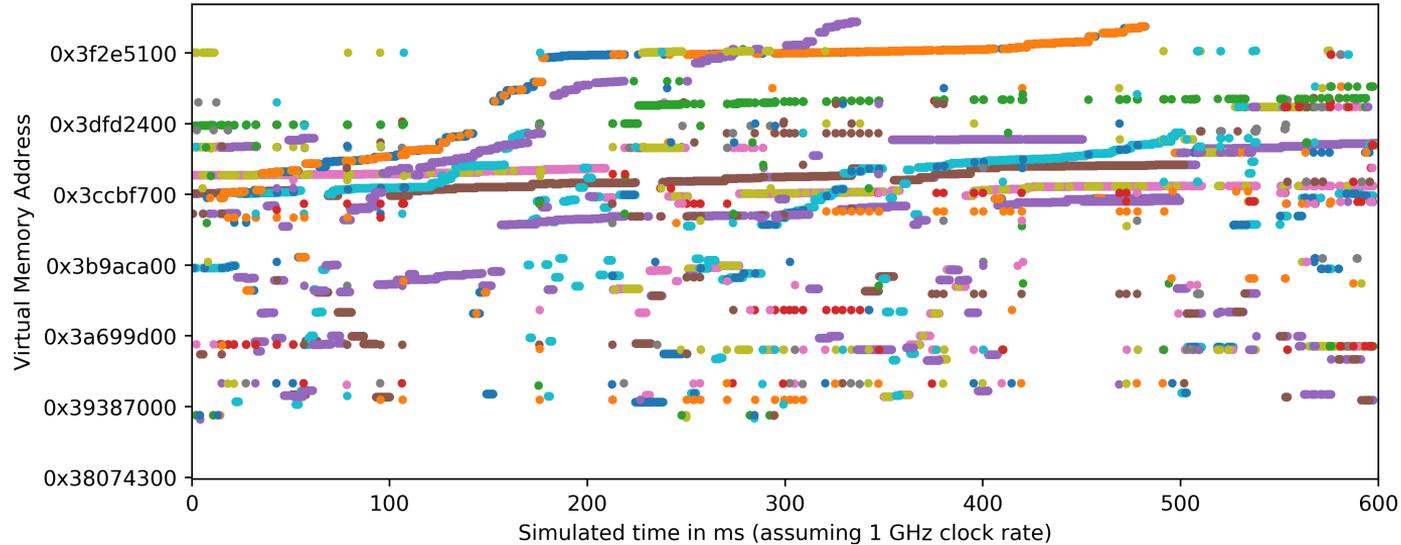


**Sampling Rate: 1 KHz**



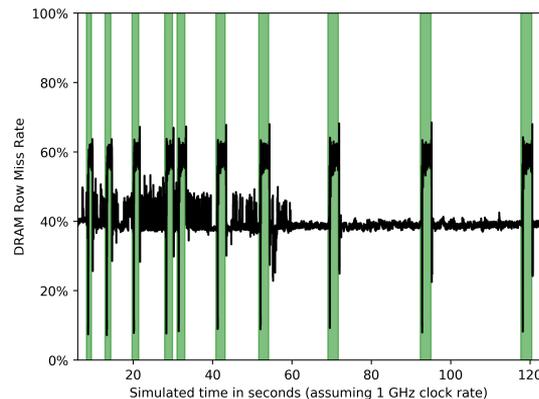
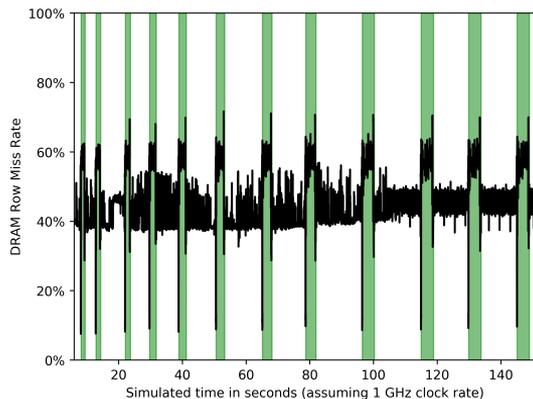
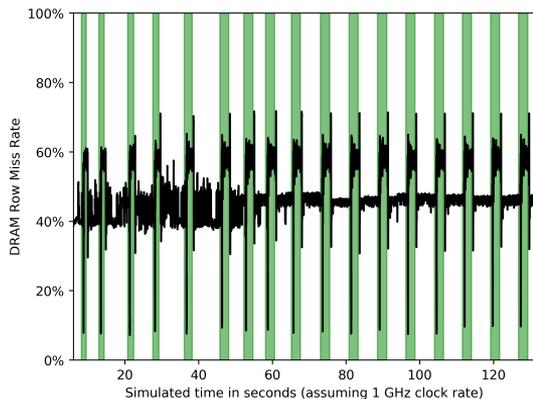
**Every Allocation**

# Logging Memory Allocations



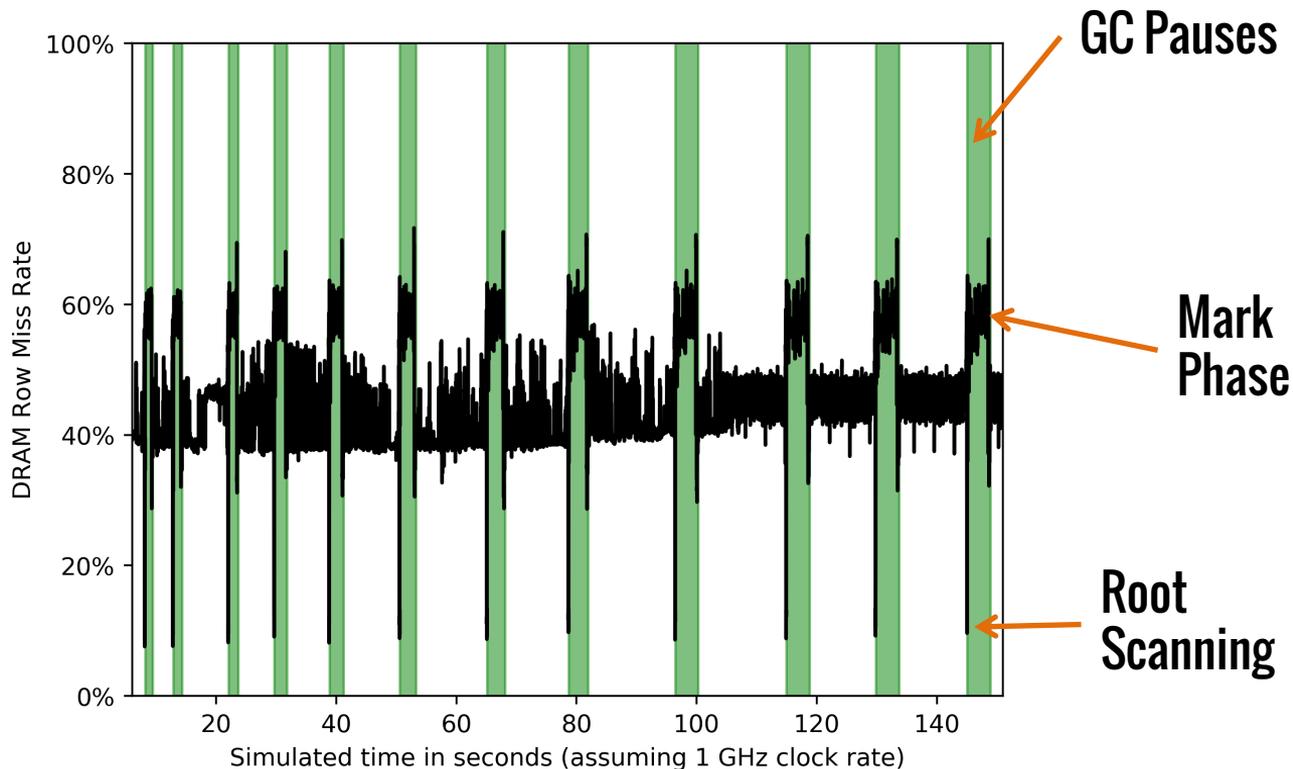
All memory allocations in a program  
(color indicates allocation class size)

# DRAM Row Misses



**Dacapo** Java Benchmarks on FPGA RISC-V core, FCFS open-page memory access scheduler. **800 Billion cycles @ 30MHz**

# DRAM Row Misses



# PART III

1. **Running JikesRVM on Rocket Chip**  
Executing JikesRVM on FPGA-based RISC-V hardware
2. **Managed-Language Use Cases**  
New research that is enabled by this infrastructure
3. **The State of Java on RISC-V**  
Progress, Challenges and Announcements

# JVM on RISC-V Progress

- **Jikes Research JVM:**
  - Baseline JIT, no optimizing JIT port yet
- **OpenJDK HotSpot JVM:**
  - Runs with zero backend, but no high-performance JIT compiler port yet



# We need your help!

Are you interested  
in working on the  
**OpenJDK** port?

# Announcement



The RISC-V Foundation is launching a new **J Extension Work Group** to add managed-language support to RISC-V!

If you would like to get involved, talk to me or David Chisnall ([david.chisnall@cl.cam.ac.uk](mailto:david.chisnall@cl.cam.ac.uk))