

# TRASH DAY: COORDINATING GARBAGE COLLECTION IN DISTRIBUTED SYSTEMS

**Martin Maas\*** † Tim Harris† Krste Asanovic\* John Kubiataowicz\*

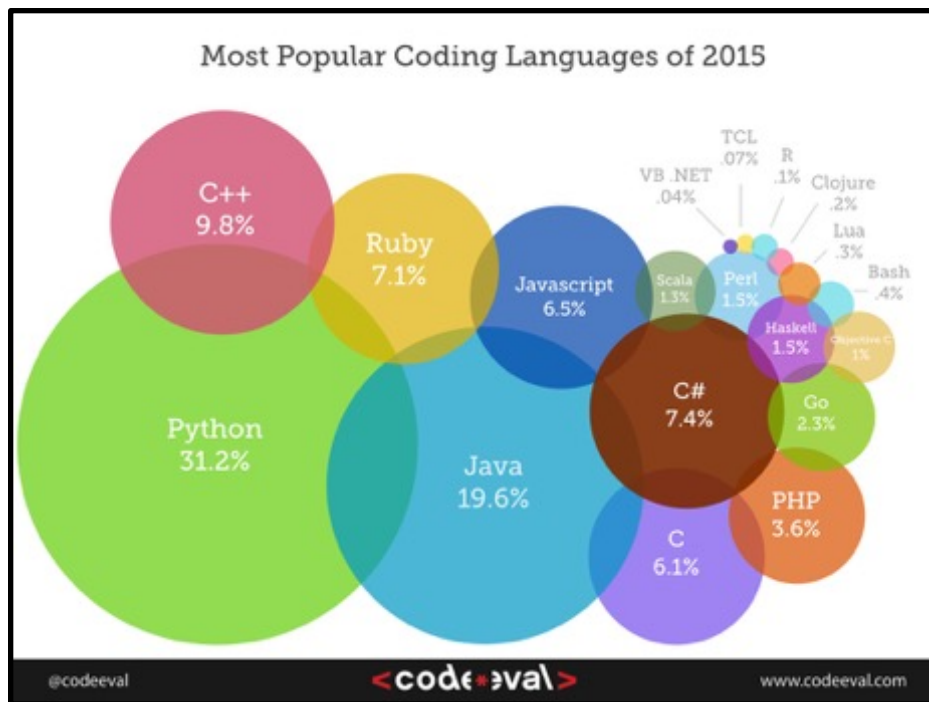
\*University of California, Berkeley † Oracle Labs, Cambridge



# Why you should care about GARBAGE COLLECTION in Data Center Applications

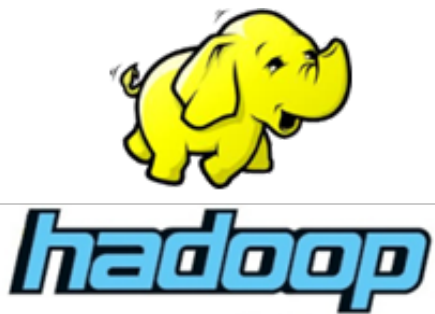


# Most Popular Languages 2015



5 out of the top 6 languages popular in 2015 use Garbage Collection (GC)

# Popular Frameworks using GC



# GC used by Cloud Companies



# Why Managed Languages?



# Productivity Gains



# Avoiding Bugs

Java bytecode	Intermediate representation 1	Intermediate representation 2
<pre>bpc 0 aload 0 1 iload 1 2 iaload 3 istore 3 4 aload 0 5 iload 2 6 iaload 7 istore 4</pre>	<pre>bpc mid (bounds check omitted for clarity...) 2 [1] const 12 ↓ [2] add [L0] [1] [3] const 2 ↓ [4] lsl [L1] [3] [5] load #1 [2]+[4] ; [6] lsl [L2] [3] ↓ [7] load #1 [2]+[6] ;</pre>	<pre>bpc mid (bounds check omitted for clarity...) 2 [1] const 2 ↓ [2] lsl [L1] [1] [3] add [L0] [2] ↓ [4] load #1 [3]+[2] ; [5] lsl [L1] [1] ↓ [6] add [L0] [2] ↓ [7] load #1 [3]+[2] ;</pre>
<p><b>Equivalent C:</b></p> <pre>L3 = L0[L1]; L4 = L0[L2];</pre>	<p><b>SPARC code</b></p> <pre>add %l0, 12, %g1 all %l1, 2, %g2 ld [%l1+%g2], %g3 all %l2, 2, %g2 ld [%l1+%g2], %g4</pre> <p><b>ARM code</b></p> <pre>add r3, r0, 12 ldr r4, [r3, r1, LSL#2] ldr r5, [r3, r2, LSL#2]</pre> <p>① ARM and SPARC support <math>[1 + r2]</math> addressing</p>	<p><b>MIPS code</b></p> <pre>lsl \$t1, \$a1, 2 add \$t0, \$t1, \$a0 lw \$a0, 12(\$t0) lsl \$t1, \$a2, 2 add \$t0, \$t1, \$a0 lw \$a1, 12(\$t0)</pre> <p>② MIPS only supports <math>[+ offset]</math> addressing rule that offset can be merged directly into the expression since these cannot change during compilation</p>

array object layout

```

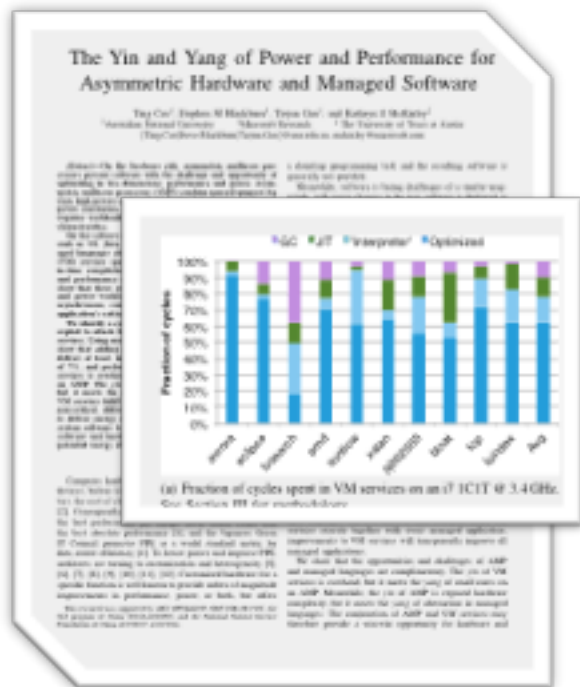
0 .S.header
4 .S.length
12 .S.data
20 .S.next
32 .S.mark

```

[Targeting Dynamic Compilation for Embedded Environments, Michael Chen and Kunle Olukotun, JMW02]

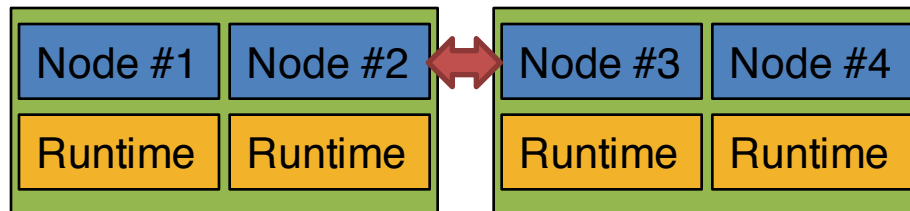
# Enable Certain Optimizations

# What is the Cost of GC?



ISCA'12: Cao et al.

- GC overhead workload and heap-size dependent, **5-20% on single machine**
- In Distributed Applications, additional overheads emerge. **Applications run across independent runtime systems:**



# Two Example Workloads



Throughput-oriented  
Batch-style

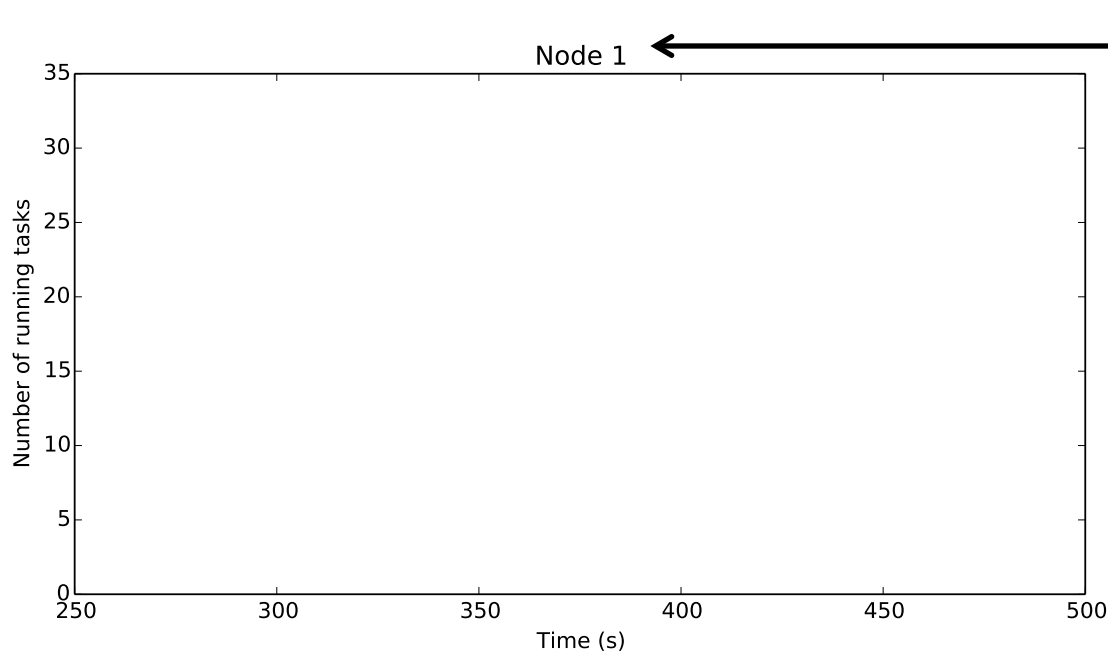


*cassandra*  
Latency-sensitive  
Interactive





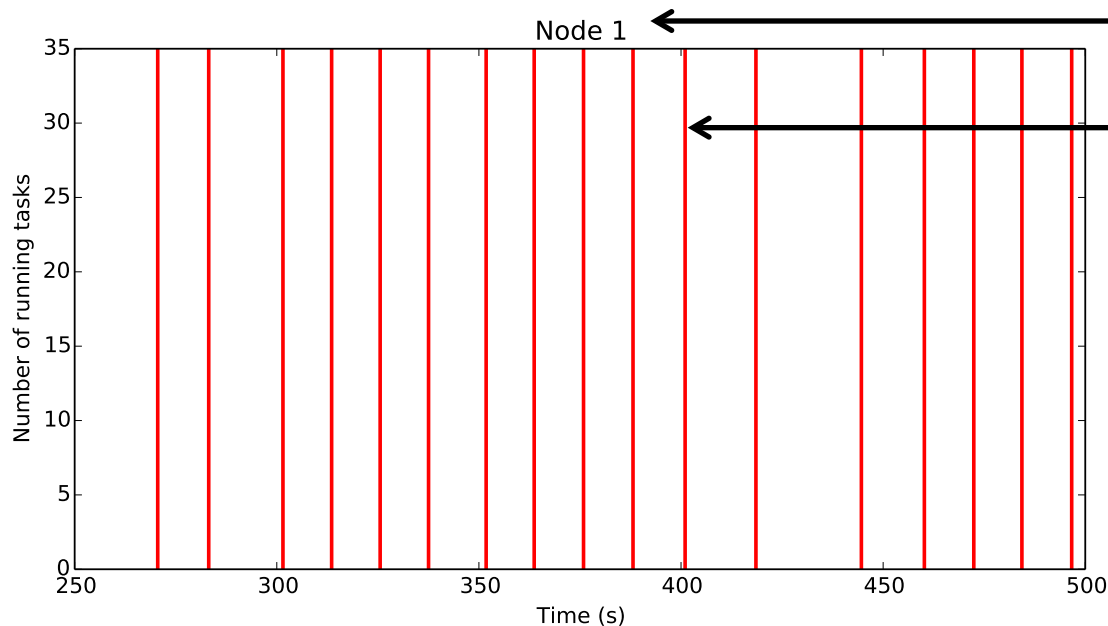
# Spark Running PageRank



8-node cluster

PageRank on 56 GB  
Wikipedia web graph

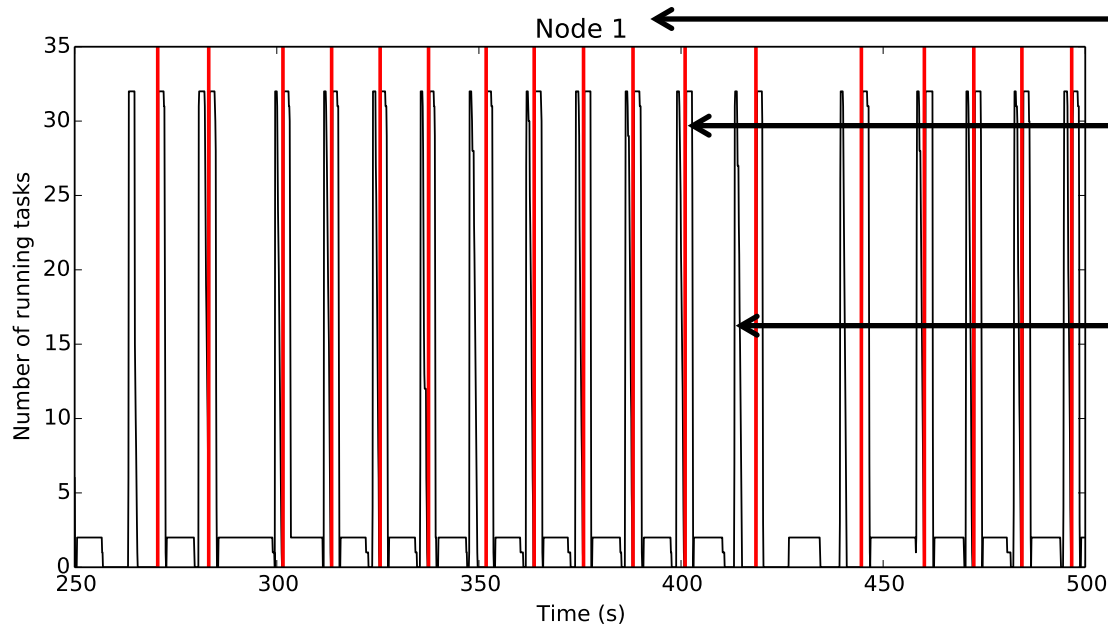
# Spark Running PageRank



8-node cluster

Execution is divided into supersteps

# Spark Running PageRank

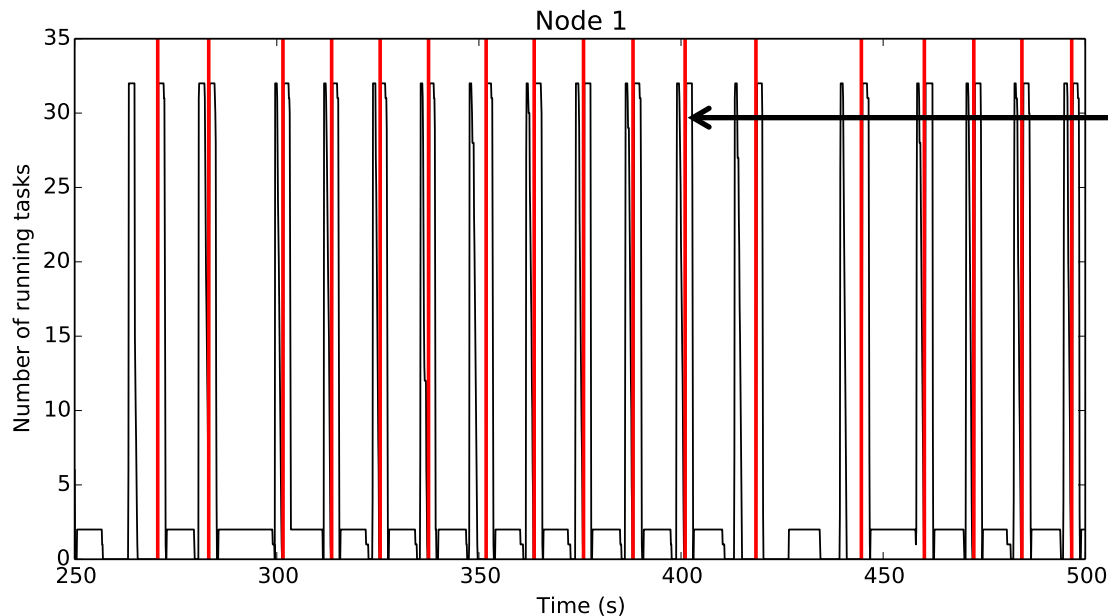


8-node cluster

Execution is divided into supersteps

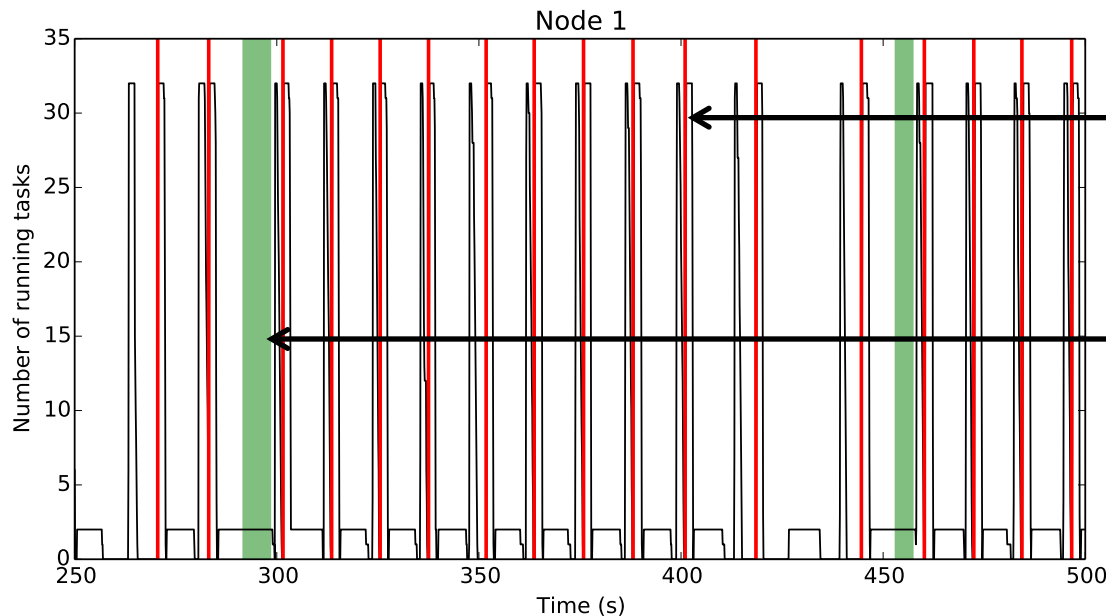
Each superstep runs independent tasks

# Spark Running PageRank



**Red** - Synchronization  
at end of superstep

# Spark Running PageRank

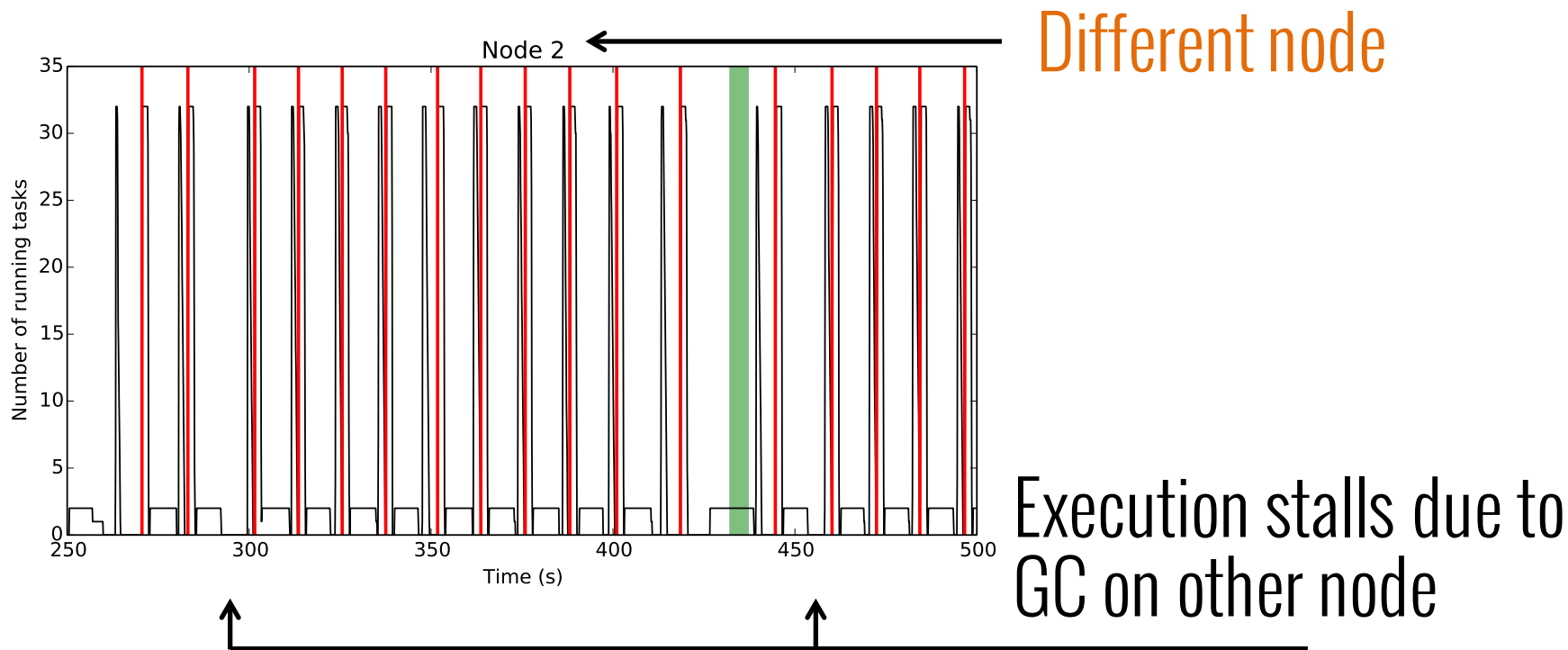


**Red** - Synchronization  
at end of superstep

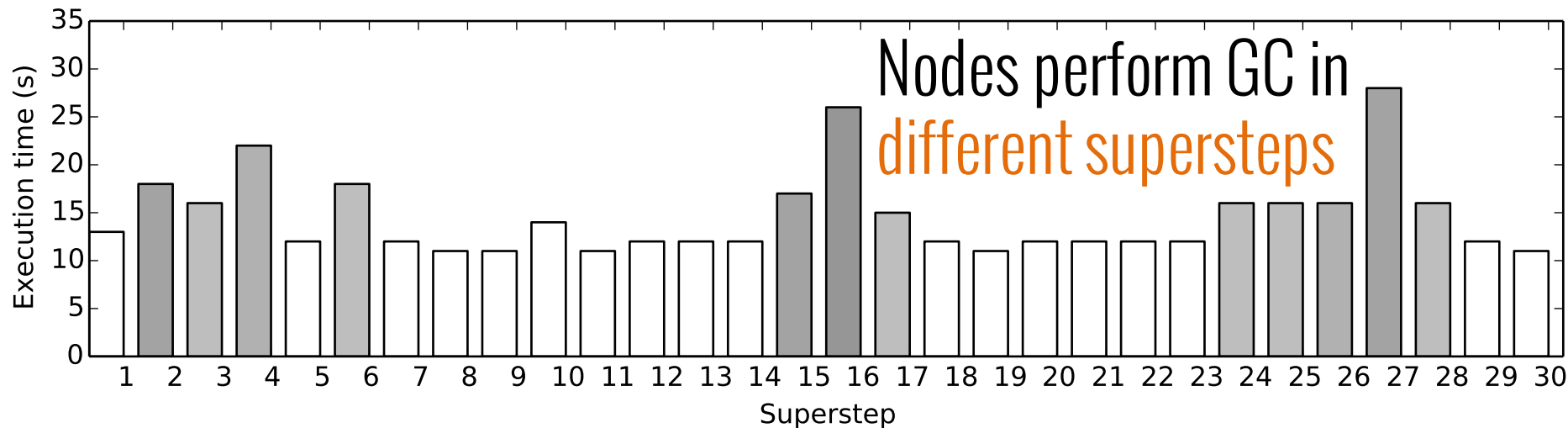
**Green** - Major GC Pause

GC prevents superstep  
from completing

# Spark Running PageRank



# Impact on Superstep Times



White = No GC during superstep

Dark = One or more GCs (the darker the more GCs)

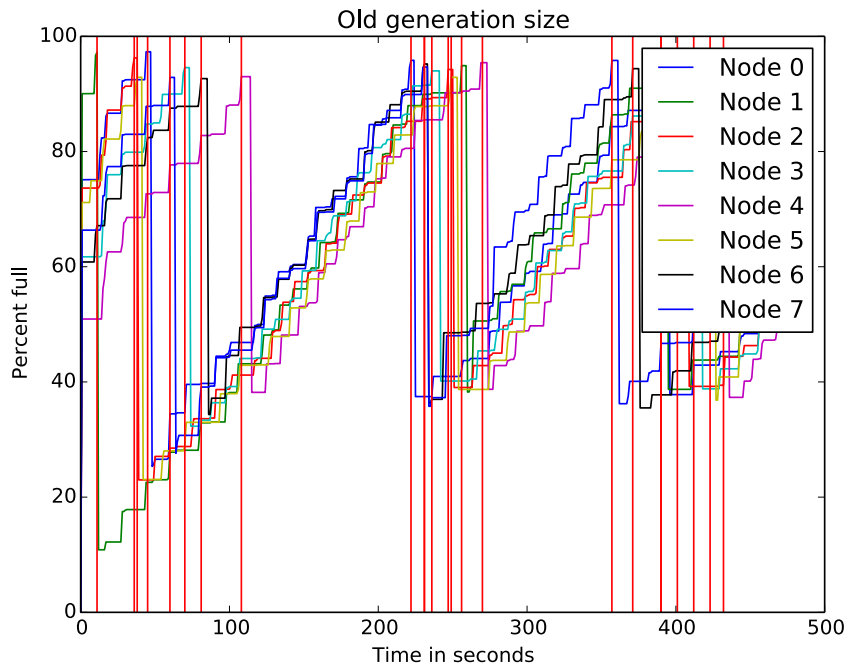


**Idea: Coordinate GC on different nodes**

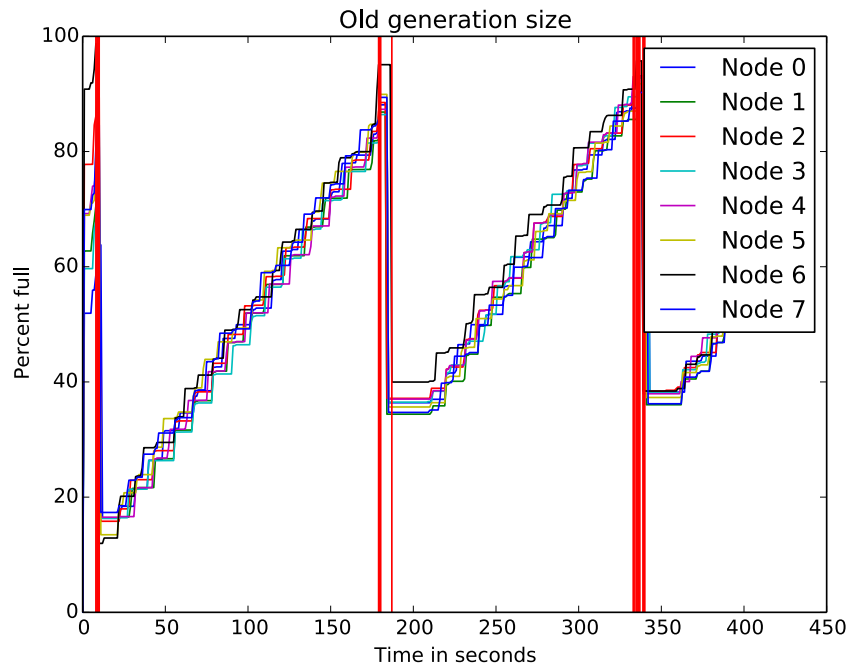
**Trigger collection on all nodes at the  
when any one reaches a threshold**

**Policy: Stop-the-world Everywhere, STWE**

# Memory Occupancy over Time

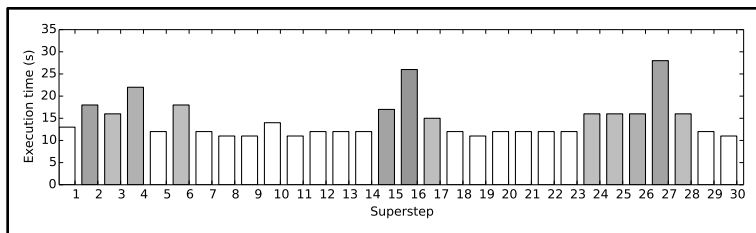
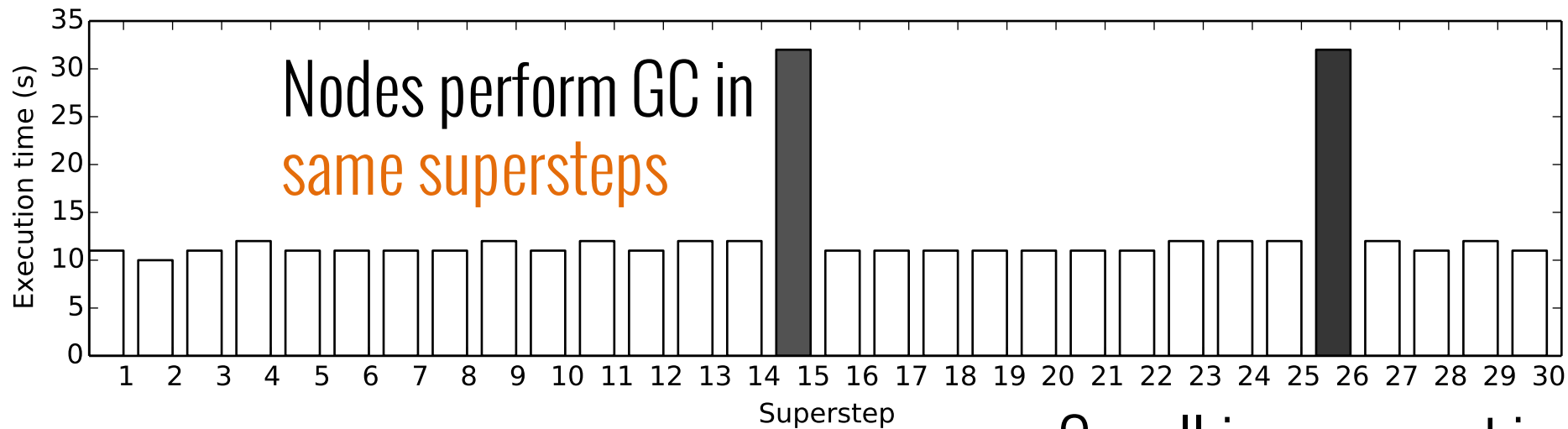


Without STWE



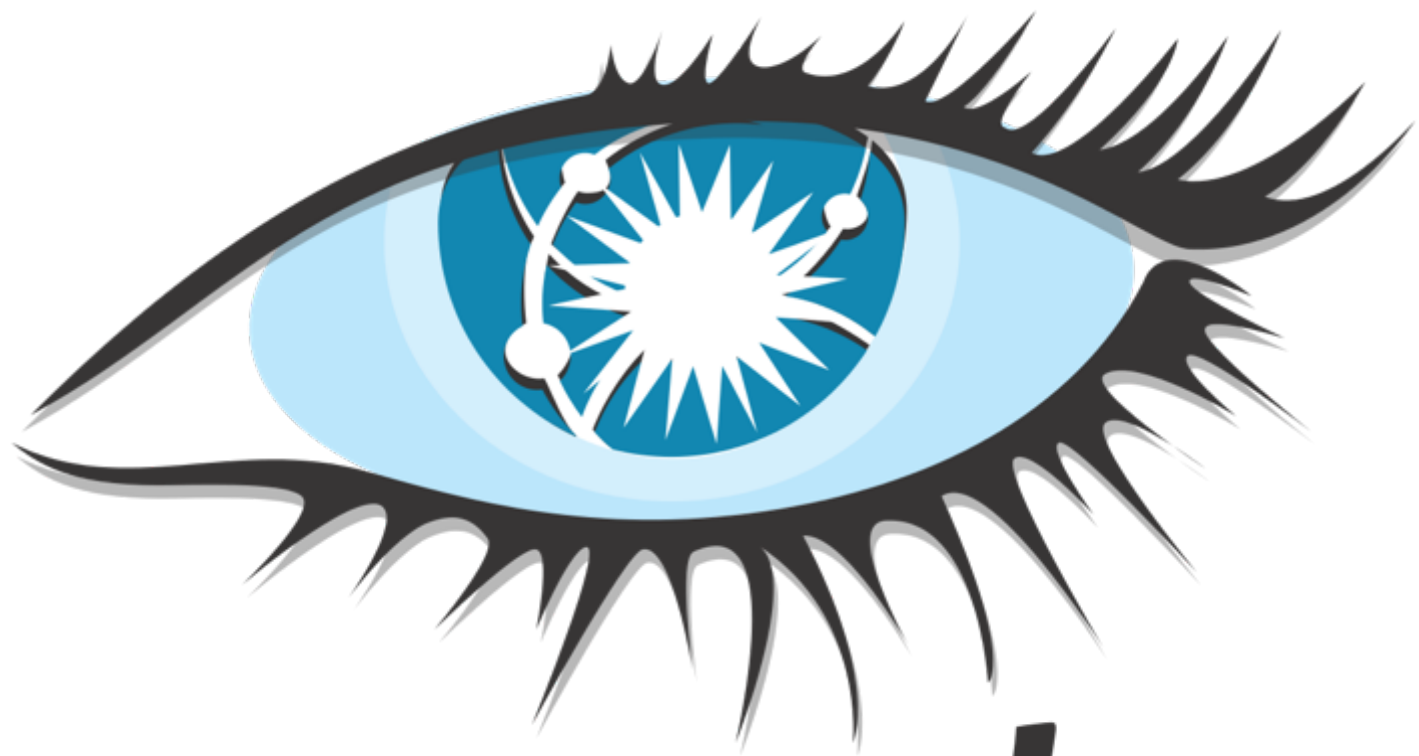
With STWE

# Impact of STWE Policy



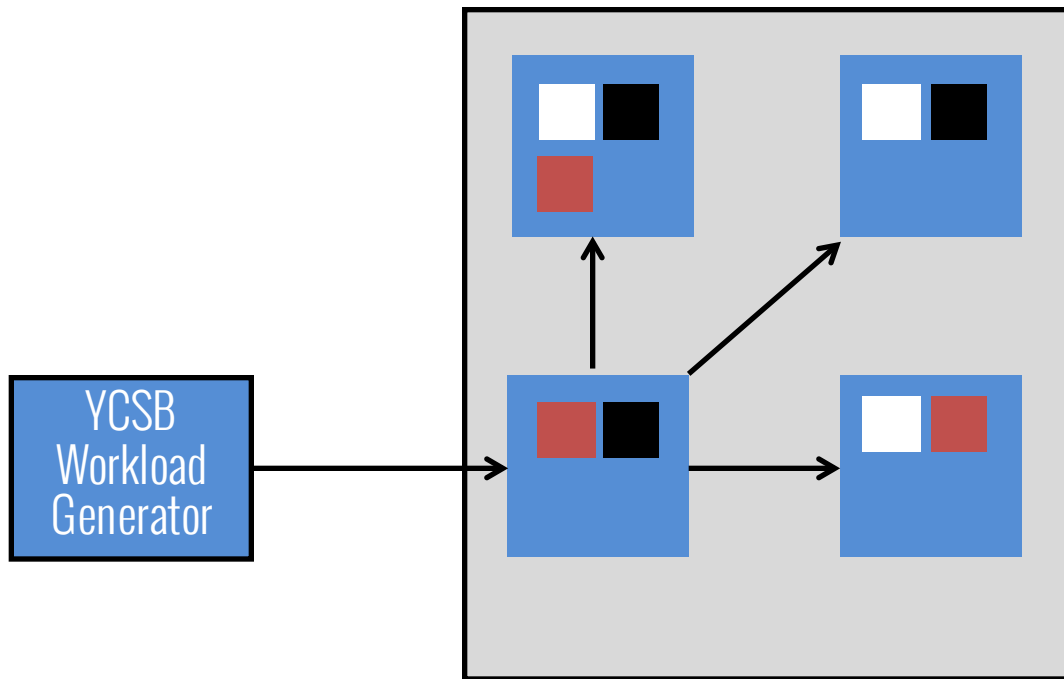
Before

Overall improvement in  
execution time (~15%)



***cassandra***

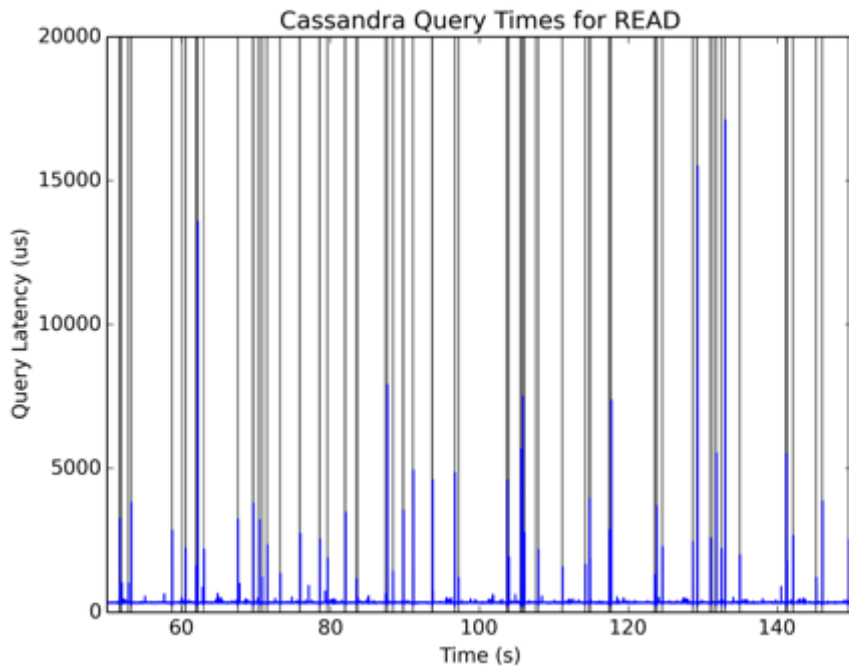
# Cassandra with YCSB



4-node Cassandra Cluster  
3-way replicated

Requests sent to arbitrary node; becomes **coordinator** and contacts replicas to **assemble quorum**.

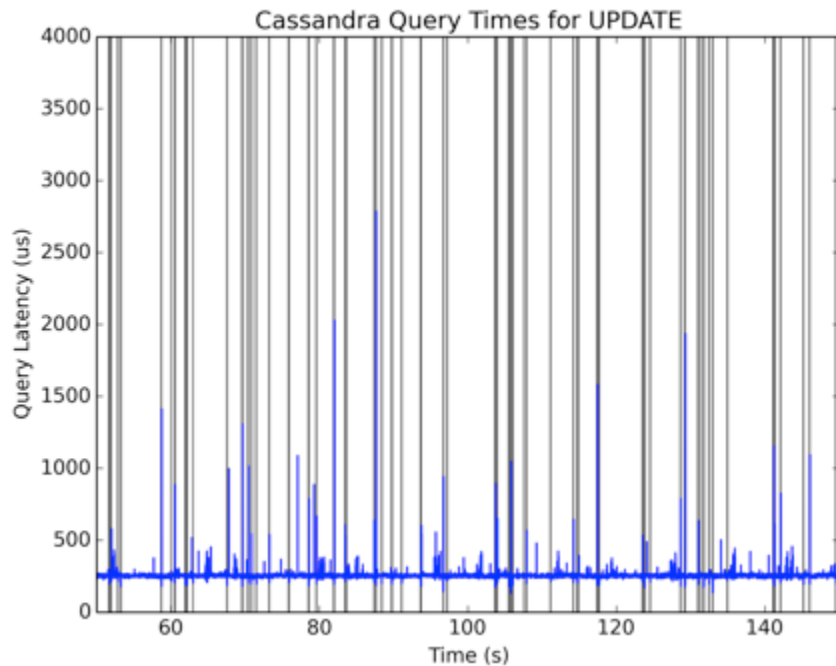
# Query Latencies over Time



Blue – mean latency  
over a 10ms window

Grey Bars – minor GC on  
any node in the cluster

# Query Latencies over Time



Blue – mean latency  
over a 10ms window

Grey Bars – minor GC on  
any node in the cluster

# Sources of Stragglers

1. Coordinator incurs GC during request
2. Node required a quorum incurs GC
3. Non-GC reasons (e.g., anti-entropy)



# Sources of Stragglers

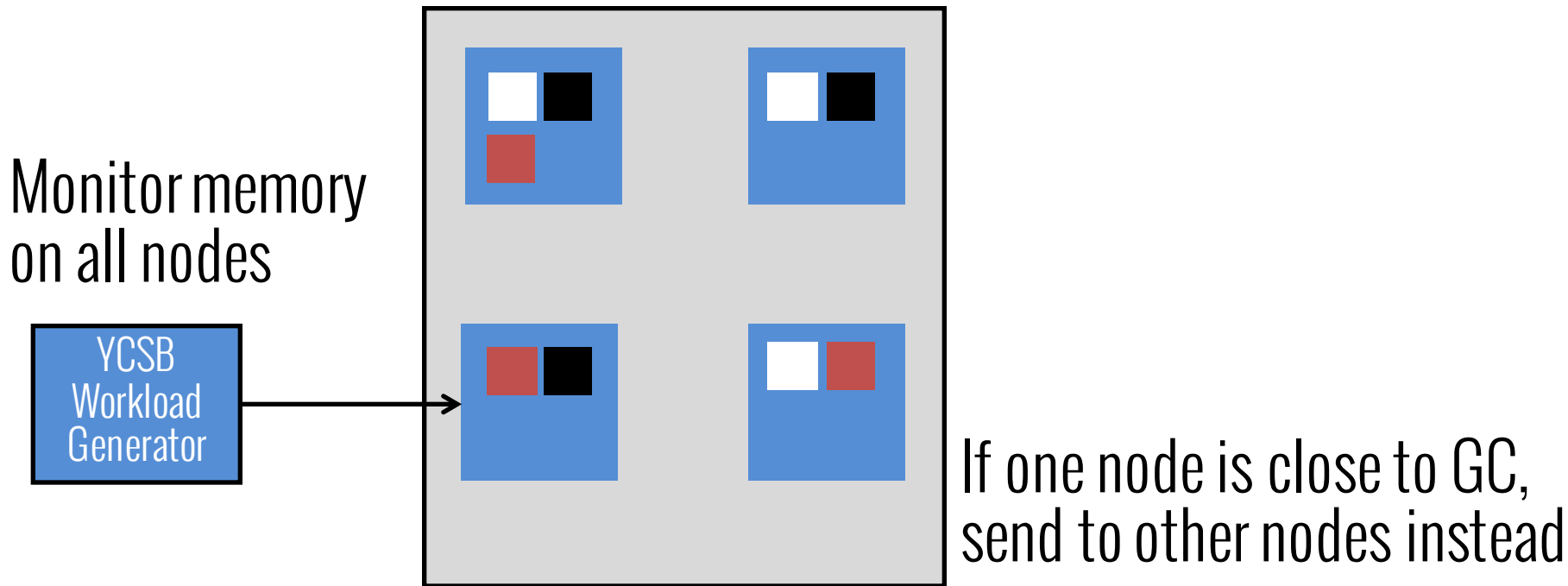
1. Coordinator incurs GC during request
2. Node required a quorum incurs GC
3. Non-GC reasons (e.g., anti-entropy)

# GC-aware Work Distribution

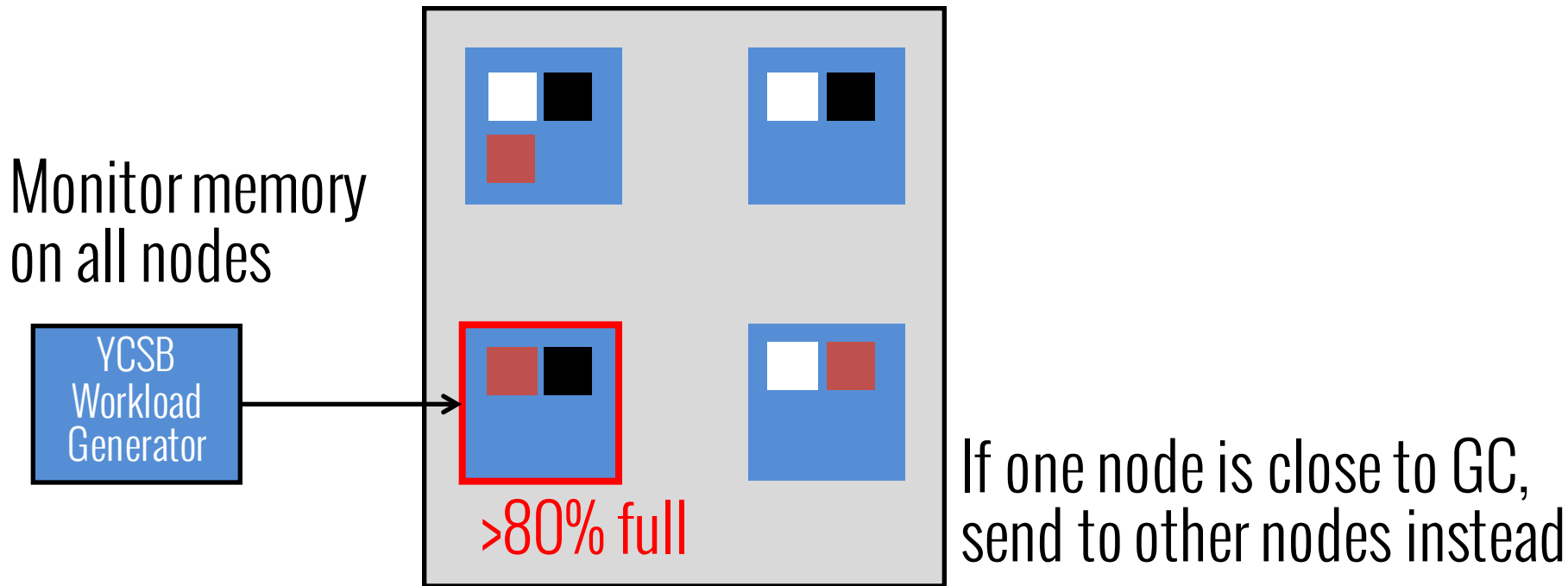
**Steer client requests to Cassandra nodes, avoiding ones that will need a minor collection soon**

Policy: Request Steering, STEER

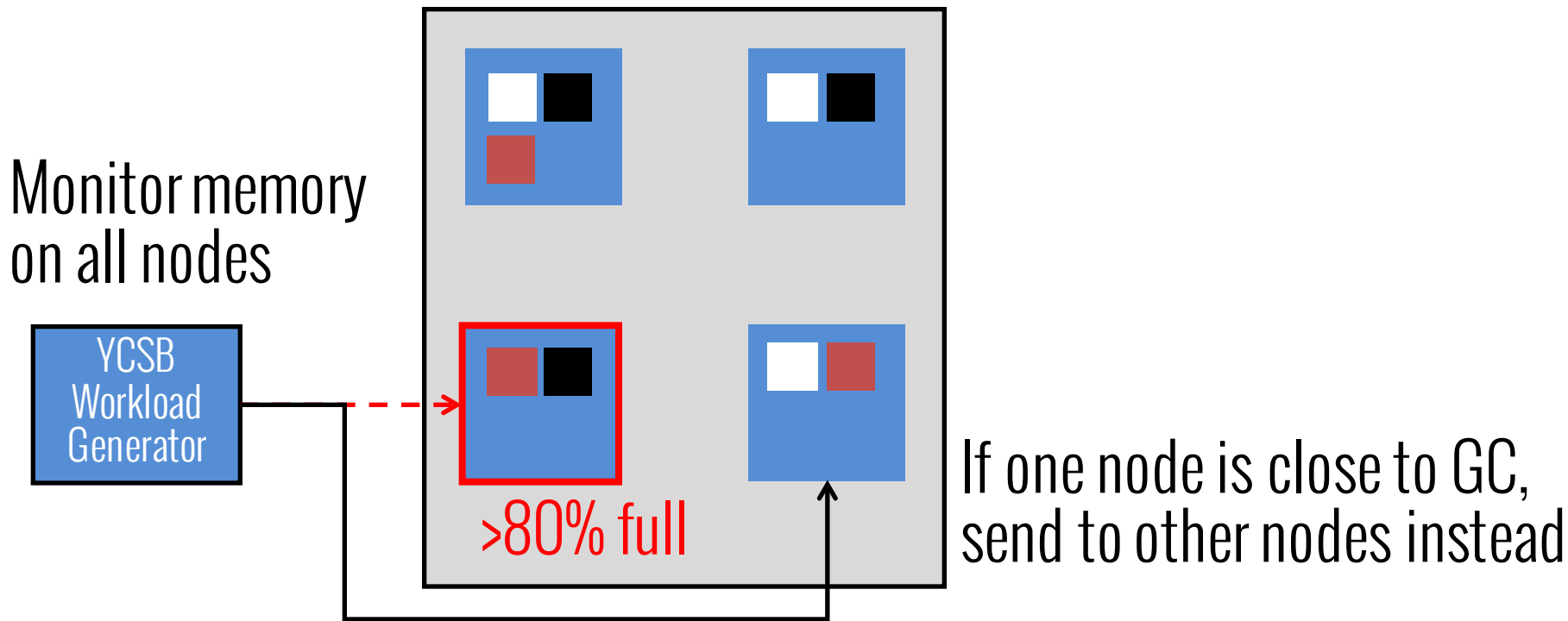
# Steering Cassandra Requests



# Steering Cassandra Requests

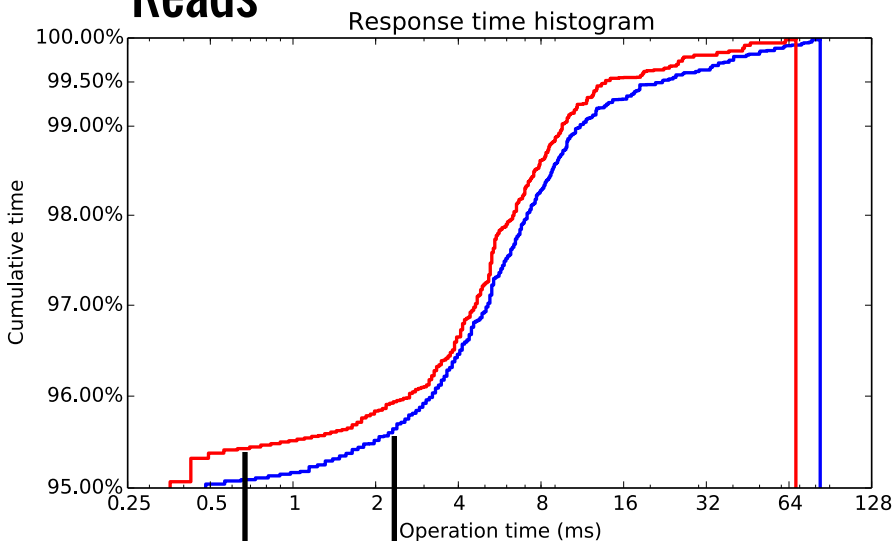


# Steering Cassandra Requests



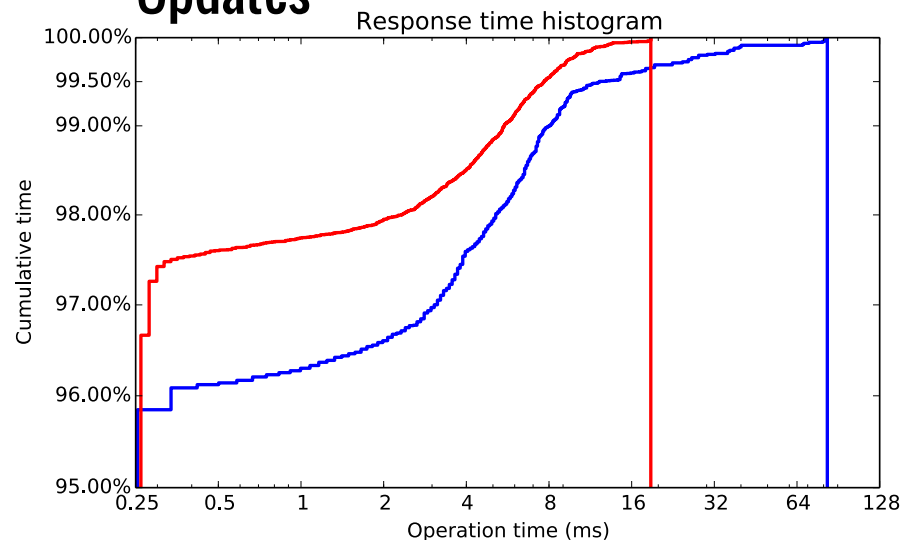
# Impact of Request Steering

## Reads



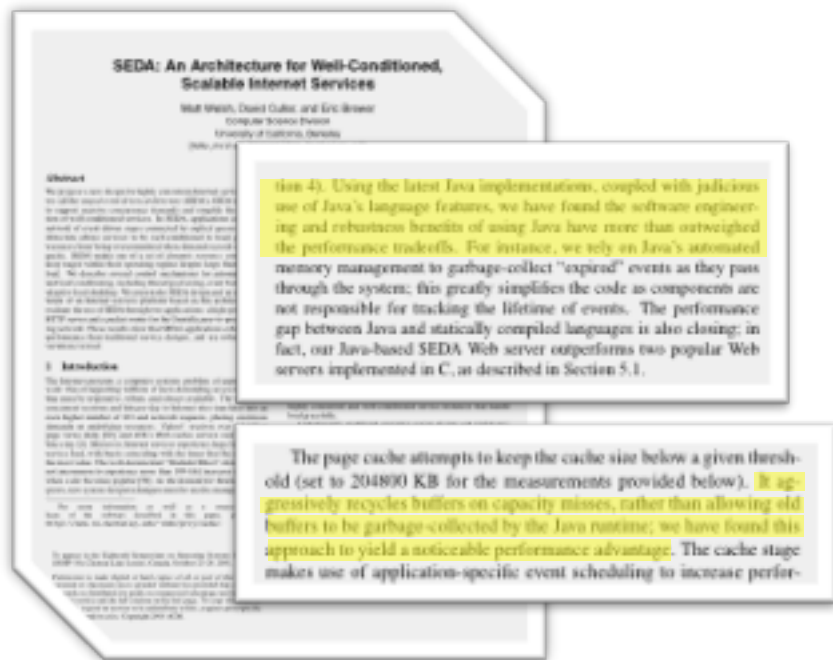
Blue – without steering  
Red – with steering

## Updates



99.9 percentile: 3.3 ms -> 1.6 ms  
Worst case: 83 ms -> 19 ms


# Are These Problems Common?



SOSP '01, Welsh et al.

- GC problems affect a large number of applications
- Have existed since dawn of warehouse-scale computing
- Current surge of interest in both industry and academia (6 new papers in last 4 mo.)


# Common Solutions



✗ Lose language advantages, lack of generality

- Rewrite in C/C++
- Use non-idiomatic language constructs

Rewrite at lower level



✗ Substantial effort to adopt

Respond to GC Pauses



✗ Performance overheads, still have pauses

Concurrent Collectors



# Common Solutions

**No general, widely  
adopted solution!**

Rewrite at  
lower level

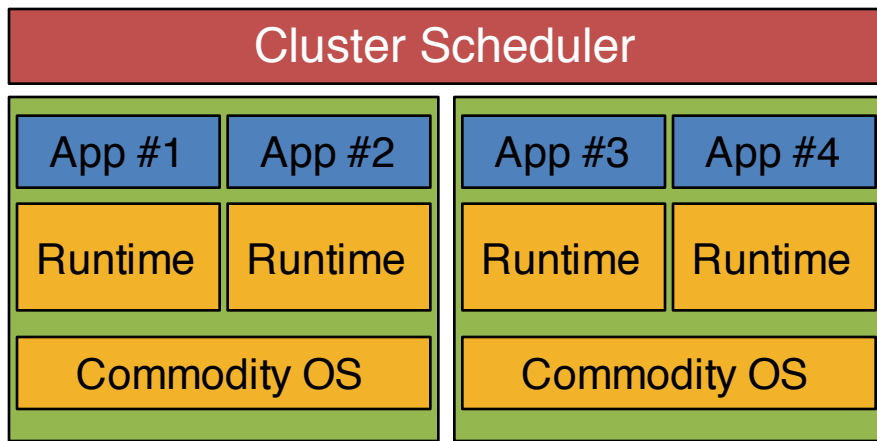
Respond to  
GC Pauses

Concurrent  
collectors

**The problem is not GC, it is  
language runtime system coordination**

# Current Approach

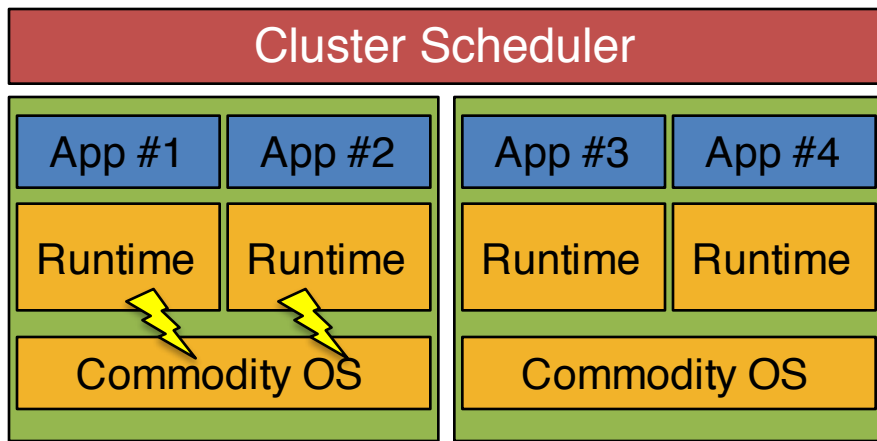
Language Runtime Systems are  
completely independent (not just GC)



# Current Approach

Language Runtime Systems are  
completely independent (not just GC)

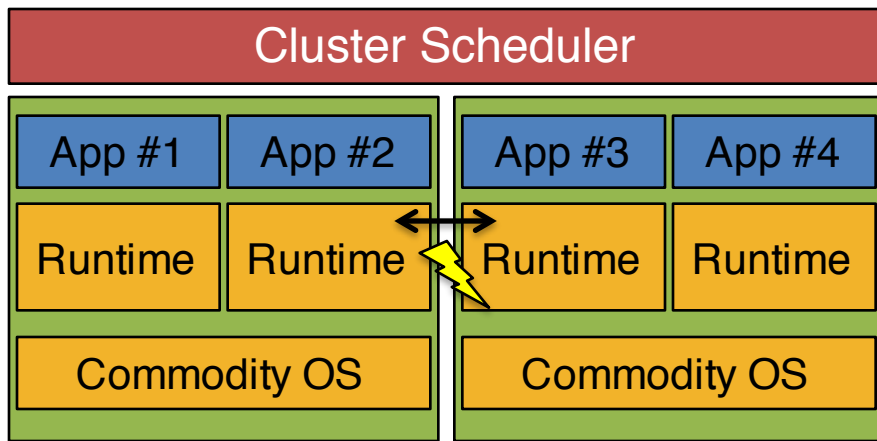
Intra-node  
Interference



# Current Approach

Language Runtime Systems are  
completely independent (not just GC)

Intra-node  
Interference

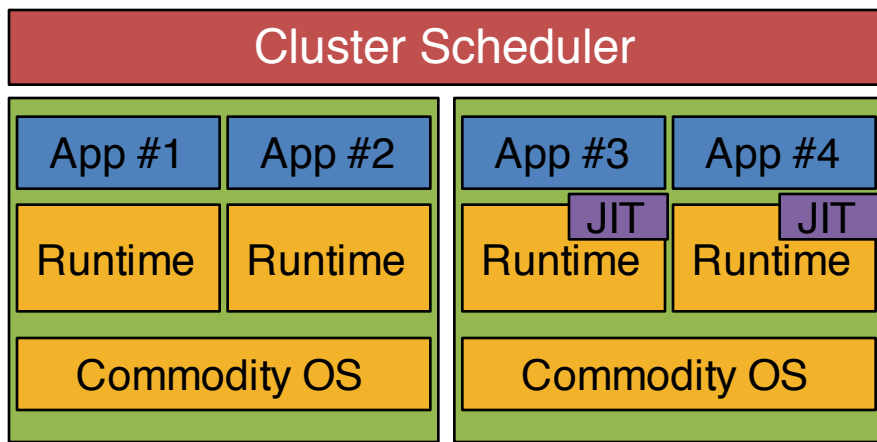


Lack of  
Coordination

# Current Approach

Language Runtime Systems are  
completely independent (not just GC)

Intra-node  
Interference



Redundancy

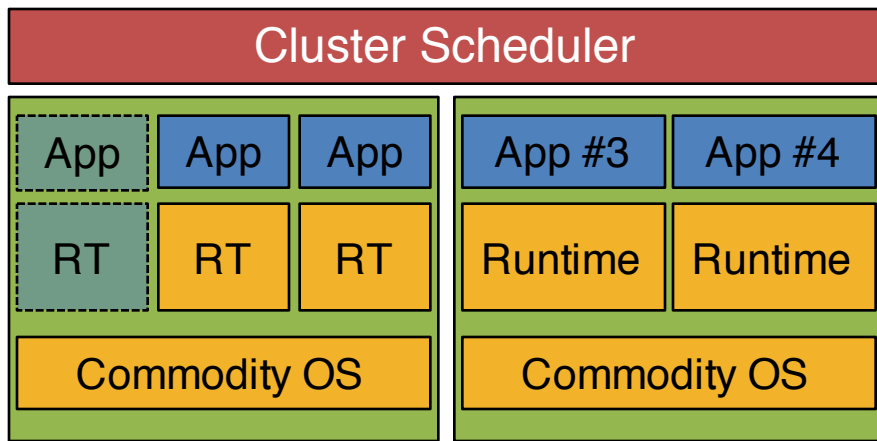
Lack of  
Coordination

# Current Approach

Language Runtime Systems are  
completely independent (not just GC)

Elasticity

Intra-node  
Interference

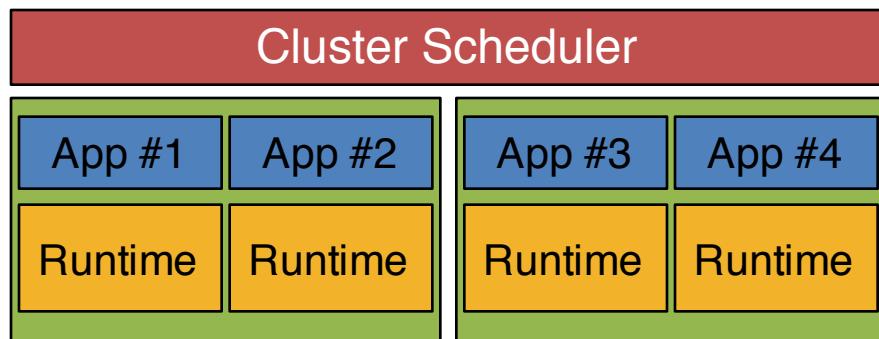


Redundancy

Lack of  
Coordination

# Holistic Runtime Systems

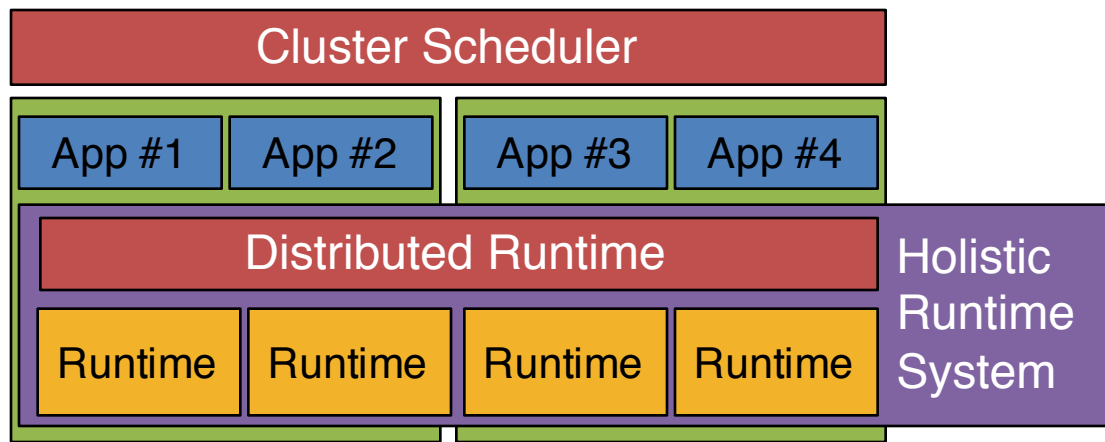
Apply the Distributed OS Ideas to design  
a **Distributed Language Runtime System**





# Holistic Runtime Systems

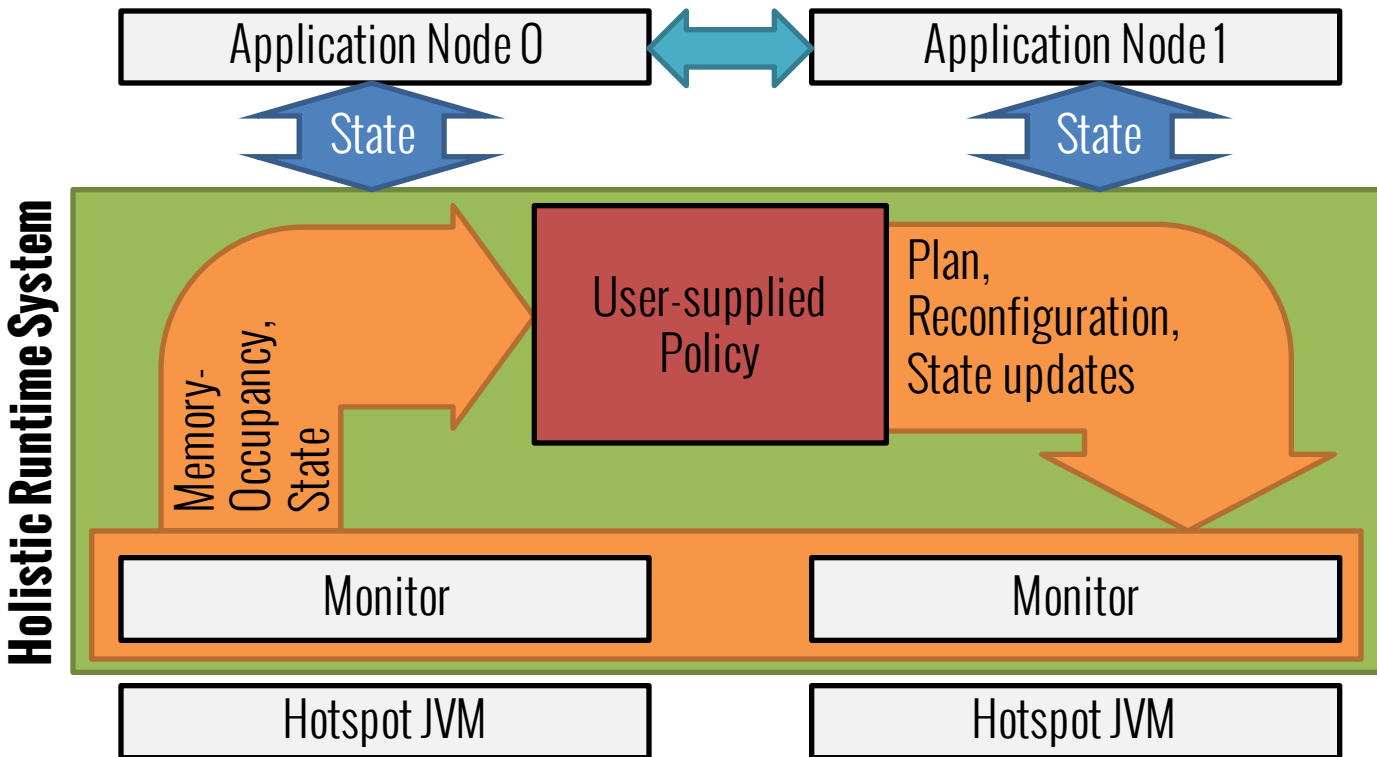
Apply the Distributed OS Ideas to design  
a **Distributed Language Runtime System**



# Our Prototype

- Coordinated runtime decisions using a **feedback loop with dist. consensus**
- Configured by **Policy** (written in DSL)
- **Drop-in replacement** for Java VM
- **No modifying of application** required

# System Design



# Why This Approach?

- **Easy to adopt** (just pick policy, almost no configuration required)
- **Minimally invasive** to runtime system
- **Expressive** (can express a large range of GC coordination policies)

**Our plan is to make the system available  
as open source**

**Would you use it?**

# Thank you! Any Questions?



**Martin Maas, Tim Harris, Krste Asanovic, John Kubitowicz**

{maas,krste,kubitron}@eecs.berkeley.edu

timothy.l.harris@oracle.com

**Work started while at Oracle Labs, Cambridge.**