Improving the Design and the Performance of Managed Runtime Environments

Gaël Thomas

INRIA/LIP6

Exploiting Multicores in Managed Runtime Environments

Lokesh Gidra, Gaël Thomas, Julien Sopena Marc Shapiro

INRIA/LIP6

Managed Runtime Environments (MREs)



MRE: simulates an abstract hardware/OS

- ✓ **Safety**: isolate code from the rest of the system
- ✓ **Portability**: write once, run anywhere

MREs are efficient



MREs are everywhere



Smartphones



Web browsers



Desktop



Web servers

But they were not prepared to multicore

Most MREs were designed for a monocore architecture ⇒ Necessary to study their bottlenecks on a multicore architecture



AMD Opteron 6172 Magny-Cours

GC Scalability (Lusearch) [PLOS'11]













The concurrency issue \Rightarrow G freed while still used



Common believe



Current believe:

STW are unacceptable for server apps [Iyengar, ISMM 2012] Long pauses due to larger heaps

Our hypothesis

Increase in transistor count is for both memory and CPU

- ✓ Large heaps come with large core count
- ✓ STW GC should be still useful, provided they scale

Can we make a GC scales with the number of cores to avoid the price of concurrent collectors?

Contribution

Identify the bottlenecks of Parallel Scavenge

(the most scalable GC of OpenJDK – used by default)

- ✓ Heavy contended locks
- ✓ Lack of NUMA-awareness

Solve the bottlenecks

- \checkmark Remove all the locks during the collection
- ✓ Propose 3 NUMA-aware heap layouts
 - Interleaved: balance memory accesses across the nodes
 - Fragmented: balance + increase memory locality
 - Segregated: balance + perfect memory locality

- 1. Background
- 2. The lock bottleneck
- 3. The NUMA bottleneck
- 4. Evaluation
- 5. Conclusion

- 1. Background
- 2. The lock bottleneck
- 3. The NUMA bottleneck
- 4. Evaluation
- 5. Conclusion



Step 1: identify the root objects (globals, stack)



Step 2: copy an object from the pending queue + update pending queue



Step 2: copy an object from the pending queue + update pending queue



Step 2: copy an object from the pending queue + update pending queue



Step 2: copy an object from the pending queue + update pending queue



Step 2: copy an object from the pending queue + update pending queue



Step 2: copy an object from the pending queue + update pending queue



Step 3: invert the spaces + consider to space empty



Advantage: spaces are never fragmented

- 1. Background
- 2. The lock bottleneck
- 3. The NUMA bottleneck
- 4. Evaluation
- 5. Conclusion

Poor Synchronization in Parallel Scavenge



Coarse grained synchronization + use of monitors

Simplify synchronizations



- 1. Background
- 2. The lock bottleneck
- 3. The NUMA bottleneck
- 4. Evaluation
- 5. Conclusion

Impact of a NUMA architecture

Problem 1: unbalanced memory accesses Interconnect or memory controllers saturate



Impact of a NUMA architecture

Problem 2: remote memory accesses

Interconnect saturates



Impact of a NUMA architecture



Inefficient Memory Layout in ParallelScavenge (PS)

The initial thread fixes the mapping of physical pages



SPECjbb2005 allocates ~95% of memory from a single node

Gaël Thomas

Solution 1: Interleaved Space

Map the pages on the node in round-robin



\Rightarrow perfect memory balance

Solution 1: Interleaved Space

Idea: map the pages on the node in round-robin



Associate fragments to memory nodes

 \Rightarrow node-local allocation



Naturally balance the load

- \checkmark GC threads uniformly distributed on the nodes
- ✓ Efficient work stealing between GC threads
- \Rightarrow balance the copies on all the nodes



Increase locality for the application

- ✓ Mutator mostly accesses recently allocated objects
- \checkmark Recently allocated object is on the mutator's node





Solution 3: Segregated Space

Fragmented space + node-local scanning (Send remote references to the owner of the object)

Perfect locality for the GC

But have to pay the price of inter-node message exchanges

Good locality for the mutators Mutator mostly accesses recently allocated objects

Natural balance of the load if allocation rates of the mutators are similar

Summary of the spaces

	Interleaved	Fragmented	Segregated
Memory Balance	Perfect	Good	Good
Mutator Locality	Bad	Good	Good
GC Locality	Bad	Good	Perfect
Comment			Inter-node messages

- 1. Background
- 2. The lock bottleneck
- 3. The NUMA bottleneck

4. Evaluation

5. Conclusion

Evaluation

Hardware:



Focus on SPECjbb2005 in the presentation















Overall effect



Scalability with the cores



Scalability with the cores



Scalability with the cores





To take away

STW GCs assumed to be inherently non-scalable is probably a mistake \Rightarrow Stop the userId CC still well suited for contemporary U/W

 \Rightarrow Stop-the-world GC still well suited for contemporary H/W

Most important NUMA effects

- ✓ Balancing memory accesses has the most important impact
- ✓ Increasing memory locality is required to scale
- ✓ Latency improvement due to locality negligible

Next step

 \checkmark Avoiding most of the messages between the nodes

[A Study of the Scalability of Stop-the-world Garbage Collectors on Multicores, ASPLOS 2013]