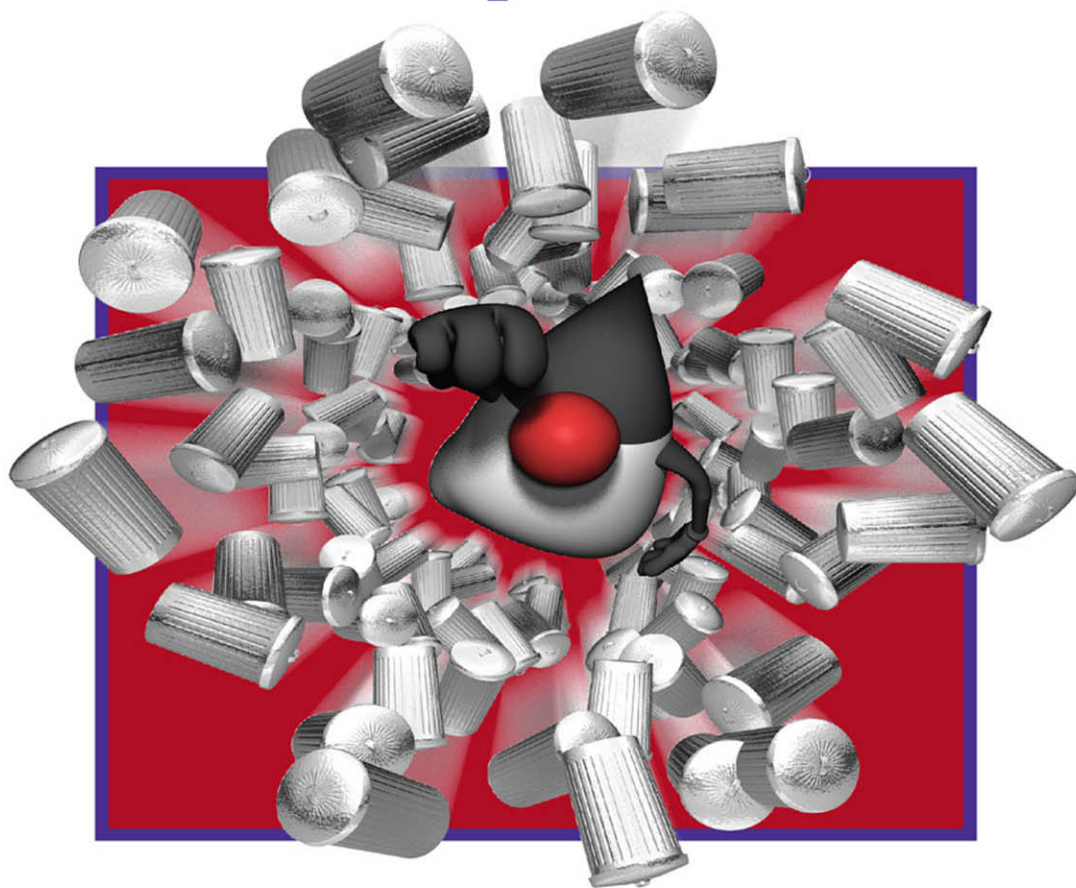


Charlie Hunt · Monica Beckwith
Poonam Parhar · Bengt Rutissson



Java[®] Performance Companion



FREE SAMPLE CHAPTER

SHARE WITH OTHERS



Java[®] Performance Companion

This page intentionally left blank



Java[®] Performance Companion

Charlie Hunt
Monica Beckwith
Poonam Parhar
Bengt Rutisson

◆ Addison-Wesley

Boston • Columbus • Indianapolis • New York • San Francisco • Amsterdam • Cape Town
Dubai • London • Madrid • Milan • Munich • Paris • Montreal • Toronto • Delhi • Mexico City
São Paulo • Sidney • Hong Kong • Seoul • Singapore • Taipei • Tokyo

Many of the designations used by manufacturers and sellers to distinguish their products are claimed as trademarks. Where those designations appear in this book, and the publisher was aware of a trademark claim, the designations have been printed with initial capital letters or in all capitals.

The authors and publisher have taken care in the preparation of this book, but make no expressed or implied warranty of any kind and assume no responsibility for errors or omissions. No liability is assumed for incidental or consequential damages in connection with or arising out of the use of the information or programs contained herein.

For information about buying this title in bulk quantities, or for special sales opportunities (which may include electronic versions; custom cover designs; and content particular to your business, training goals, marketing focus, or branding interests), please contact our corporate sales department at corpsales@pearsoned.com or (800) 382-3419.

For government sales inquiries, please contact governmentsales@pearsoned.com.

For questions about sales outside the United States, please contact intlcs@pearson.com.

Visit us on the Web: informit.com/aw

Cataloging-in-Publication Data is on file with the Library of Congress.

Copyright © 2016 Pearson Education, Inc.

All rights reserved. Printed in the United States of America. This publication is protected by copyright, and permission must be obtained from the publisher prior to any prohibited reproduction, storage in a retrieval system, or transmission in any form or by any means, electronic, mechanical, photocopying, recording, or likewise. For information regarding permissions, request forms and the appropriate contacts within the Pearson Education Global Rights & Permissions Department, please visit www.pearsoned.com/permissions.

ISBN-13: 978-0-13-379682-7

ISBN-10: 0-13-379682-5

Text printed in the United States on recycled paper at RR Donnelley in Crawfordsville, Indiana.
First printing, April 2016



Contents

Preface	ix
Acknowledgments	xi
About the Authors	xv
Chapter 1 Garbage First Overview	1
Terminology	1
Parallel GC	2
Serial GC	4
Concurrent Mark Sweep (CMS) GC	5
Summary of the Collectors	7
Garbage First (G1) GC	8
G1 Design	10
Humongous Objects	12
Full Garbage Collections	12
Concurrent Cycle	13
Heap Sizing	14
References	14

Chapter 2	Garbage First Garbage Collector in Depth	15
	Background	15
	Garbage Collection in G1	16
	The Young Generation	17
	A Young Collection Pause	18
	Object Aging and the Old Generation	19
	Humongous Regions	19
	A Mixed Collection Pause	22
	Collection Sets and Their Importance	24
	Remembered Sets and Their Importance	24
	Concurrent Refinement Threads and Barriers	28
	Concurrent Marking in G1 GC	30
	Stages of Concurrent Marking	34
	Initial Mark	34
	Root Region Scanning	34
	Concurrent Marking	34
	Remark	36
	Cleanup	36
	Evacuation Failures and Full Collection	37
	References	38
Chapter 3	Garbage First Garbage Collector Performance Tuning	39
	The Stages of a Young Collection	39
	Start of All Parallel Activities	41
	External Root Regions	42
	Remembered Sets and Processed Buffers	42
	Summarizing Remembered Sets	44
	Evacuation and Reclamation	47
	Termination	47
	Parallel Activity Outside of GC	48
	Summarizing All Parallel Activities	48
	Start of All Serial Activities	48
	Other Serial Activities	49
	Young Generation Tunables	50

	Concurrent Marking Phase Tunables	52
	A Refresher on the Mixed Garbage Collection Phase	54
	The Taming of a Mixed Garbage Collection Phase	56
	Avoiding Evacuation Failures	59
	Reference Processing	60
	Observing Reference Processing	60
	Reference Processing Tuning	62
	References	65
Chapter 4	The Serviceability Agent	67
	What Is the Serviceability Agent?	68
	Why Do We Need the SA?	68
	SA Components	69
	SA Binaries in the JDK	69
	JDK Versions with Complete SA Binaries	69
	How the SA Understands HotSpot VM Data Structures	70
	SA Version Matching	71
	The Serviceability Agent Debugging Tools	72
	HSDB	72
	HSDB Tools	80
	CLHSDB	100
	Some Other Tools	103
	Core Dump or Crash Dump Files	108
	Debugging Transported Core Files	109
	Shared Library Problems with the SA	109
	Eliminate Shared Library Problems	110
	System Properties for the Serviceability Agent	111
	Environment Variables for the Serviceability Agent	112
	JDI Implementation	113
	Extending Serviceability Agent Tools	115
	Serviceability Agent Plugin for VisualVM	117
	How to Install the SA-Plugin in VisualVM	118
	How to Use the SA-Plugin	118
	SA-Plugin Utilities	119

Troubleshooting Problems Using the SA	123
Diagnosing <code>OutOfMemoryError</code>	123
Diagnosing a Java-Level Deadlock	131
Postmortem Analysis of a HotSpot VM Crash	136
Appendix Additional HotSpot VM Command-Line Options of Interest	145
Index	155



Preface

Welcome to the *Java® Performance Companion*. This book offers companion material to *Java™ Performance* [1], which was first published in September 2011. Although the additional topics covered in this book are not as broad as the material in *Java™ Performance*, they go into enormous detail. The topics covered in this book are the G1 garbage collector, also known as the Garbage First garbage collector, and the Java HotSpot VM Serviceability Agent. There is also an appendix that covers additional HotSpot VM command-line options of interest that were not included in the *Java™ Performance* appendix on HotSpot VM command-line options.

If you are currently using Java 8, have interest in migrating to Java 8, or have plans for using Java 9, you will likely be either evaluating G1 GC or already using it. Hence, the information in this book will be useful to you. If you have interest in diagnosing unexpected HotSpot VM failures, or in learning more about the details of a modern Java Virtual Machine, this book's content on the HotSpot VM Serviceability Agent should be of value to you, too. The HotSpot VM Serviceability Agent is the tool of choice for not only HotSpot VM developers but also the Oracle support engineers whose daily job involves diagnosing and troubleshooting unexpected HotSpot VM behavior.

This book begins with an overview of the G1 garbage collector by offering some context around why G1 was implemented and included in HotSpot VM as a GC. It then goes on to offer an overview of how the G1 garbage collector works. This chapter is followed by two additional chapters on G1. The first is an in-depth description of the internals of G1. If you already have a good understanding of how the G1 garbage

collector works, and either have a need to further fine-tune G1 or want to know more about its inner workings, this chapter would be a great place to start. The third chapter on G1 is all about fine-tuning G1 for your application. One of the main design points for G1 was to simplify the tuning required to realize good performance. For instance, the major inputs into G1 are the initial and maximum Java heap size it can use, and a maximum GC pause time you are willing to tolerate. From there G1 will attempt to adaptively adjust to meet those inputs while it executes your application. In circumstances where you would like to achieve better performance, or you would like to do some additional tuning on G1, this chapter has the information you are looking for.

The remaining chapter is dedicated entirely to the HotSpot VM Serviceability Agent. This chapter provides an in-depth description of and instructions for how to use the Serviceability Agent. If you have interest in learning more about the internals of the HotSpot VM, or how to troubleshoot and diagnose unexpected HotSpot VM issues, this is a good chapter for you. In this chapter you will learn how to use the HotSpot VM Serviceability Agent to observe and analyze HotSpot VM behavior in a variety of ways through examples and illustrations.

Last, there is an appendix that includes HotSpot VM command-line options that were not included in *Java™ Performance*'s appendix on HotSpot VM command-line options. Many of the HotSpot VM command-line options found in the appendix are related to G1. And, rather than merely listing these options with only a description, an attempt is made to also mention when it is appropriate to use them.

References

[1] Charlie Hunt and Binu John. *Java™ Performance*. Addison-Wesley, Upper Saddle River, NJ, 2012. ISBN 978-0-13-714252-1.

Register your copy of *Java® Performance Companion* at informit.com for convenient access to downloads, updates, and corrections as they become available. To start the registration process, go to informit.com/register and log in or create an account. Enter the product ISBN (9780133796827) and click Submit. Once the process is complete, you will find any available bonus content under “Registered Products.”



Acknowledgments

Charlie Hunt

For those who have ever considered writing a book, or are curious about the effort involved in doing so, the book-writing experience is a major undertaking! For me it just would not have happened without the help of so many people. I cannot begin to mention everyone who made this possible.

In an attempt to at least name those who have had a profound impact on this book getting drafted and eventually into print, I would first like to thank my coauthors, Monica Beckwith, Bengt Rutisson, and Poonam Parhar. When the idea of doing a companion book to *Java™ Performance* first surfaced, I thought it would be great to offer the opportunity to these talented HotSpot VM engineers to showcase their expertise. I am sure I have learned much more from each of them than they have learned from me. I could not be prouder of their contributions to this book.

I also extend sincere thanks to Monica Beckwith for her persistence and passion in sharing her in-depth knowledge of G1 GC. In the early days of G1, I had the pleasure of working with Monica on a daily basis on G1 performance, eventually handing off full reins to her. She has done an exceptional job with driving G1's performance and sharing her G1 knowledge.

I also have to explicitly call out Poonam Parhar and thank her for her patience. Poonam so patiently waited for the other contributors to complete their initial drafts—patiently as in years of patience! Had all of us finished our drafts in a timely way, this book probably would have been on the shelf at least two years earlier.

I also extend my thanks to the entire HotSpot VM team, the HotSpot GC engineering team, and in particular the G1 GC engineers, both past and present.

And to the reviewers of the material in this book: Paul Hohensee for his relentless attention to detail and incredible suggestions for improved readability, and Tony Printezis for his thorough review of the gory details of G1 GC and recommended tuning of G1.

Thanks also to John Cuthbertson for sharing his knowledge of G1. John also happens to be one of the most talented concurrency troubleshooting engineers I have ever worked with. I don't think I ever saw a situation where I was able to stump him with some bizarre observation about G1 that was clearly some kind of concurrency bug. He was always able to track it down.

And to Bernard Traversat and Georges Saab for their support and encouragement in pulling together material for a follow-on to *Java™ Performance*.

And obviously thanks to Greg Doench, our editor, for his patience with our many delays in delivering drafts, completing reviews, and getting the manuscript in shape to put in his hands.

Last, thanks to my wife, Barb, and son, Boyd, for putting up with yet another round of the book-writing experience!

Monica Beckwith

I felt honored when I was approached by my mentor Charlie Hunt to write a few chapters for this book. I didn't have the slightest idea that it would take me so long. So, my first set of thanks goes to my fellow writers for their patience and to Charlie for his persistence and encouragement throughout. While we are talking about encouragement, I want to thank my hubby, Ben Beckwith—when he saw my frustration he had nothing but words of encouragement for me. He was also the initial reviewer of my drafts. Thank you, Ben. And then, of course, my two kiddos, Annika and Bodin, and my mom, Usha, who have been nothing but supportive of me and of this book.

My technical strength on G1 taps off John Cuthbertson, and I am thankful to him for supporting my crazy queries and patiently listening and working with me to “make G1 adaptive” and to “tame mixed collections.” When we used to discuss the adaptive marking threshold, I got tired of typing and talking about `InitiatingHeapOccupancyPercent`, so I shortened it to IHOP and John just loved it. It's really hard to find such supportive colleagues as John and Charlie.

And then there are Paul Hohensee and Tony Printezis. They are my mentors in their own right, and I can assure you that their persistence in reviewing my chapters has improved the readability and content by at least 75 percent! :)

Thank you all for trusting me and encouraging me. I am forever in your debt!

Poonam Parhar

I was deeply honored and excited when Charlie suggested that I write a chapter on the Serviceability Agent. I thought it was a great idea, as this wonderful tool is little known to the world, and it would be great to talk about its usefulness and capabilities. But I had never written a book before, and I was nervous. Big thanks to Charlie for his trust in me, for his encouragement, and for guiding me throughout writing the chapter on the SA.

I would like to thank my manager, Mattis Castegren, for always being supportive and encouraging of my work on this book, and for being the first reviewer of the chapter on the SA. Huge thanks to Kevin Walls for reviewing my chapter and helping me improve the quality of the content.

Special thanks to my husband, Onkar, who is my best friend, too, for being supportive and always being there whenever I need help. And of course I am grateful to my two little angels, Amanvir and Karanvir, who are my continuous source of motivation and happiness.

And my most sincere thanks to my father, Subhash C. Bajaj, for his infectious cheerfulness and for being a source of light, and for always inspiring me to never give up.

Bengt Rutisson

When Charlie asked me to write a chapter for this book, I was very honored and flattered. I had never written a book before and clearly had no idea how much work it is—even to write just one chapter! I am very grateful for all the support from Charlie and the reviewers. Without their help, I would not have been able to complete this chapter.

A big thanks to my wife, Sara Fritzell, who encouraged me throughout the work and helped me set up deadlines to get the chapter completed. And, of course, many thanks to our children, Max, Elsa, Teo, Emil, and Lina, for putting up with me during the writing period.

I would also like to thank all of the members of the HotSpot GC engineering team, both past and present. They are by far the most talented bunch of engineers I have ever worked with. I have learned so much from all of them, and they have all inspired me in so many ways.

This page intentionally left blank



About the Authors

Charlie Hunt (Chicago, IL) is currently a JVM Engineer at Oracle leading a variety of Java SE and HotSpot VM projects whose primary focus is reducing memory footprint while maintaining throughput and latency. He is also the lead author of *Java™ Performance* (Addison-Wesley, 2012). He is a regular presenter at the JavaOne Conference where he has been recognized as a Java Rock Star. He has also been a speaker at other well-known conferences, including QCon, Velocity, GoTo, and Dreamforce. Prior to leading a variety of Java SE and HotSpot VM projects for Oracle, Charlie worked in several different performance positions, including Performance Engineering Architect at Salesforce.com and HotSpot VM Performance Architect at Oracle and Sun Microsystems. He wrote his first Java application in 1998, joined Sun Microsystems in 1999 as Senior Java Architect, and has had a passion for Java and JVM performance ever since.

Monica Beckwith is an Independent Performance Consultant optimizing customer applications for server-class systems running the Java Virtual Machine. Her past experiences include working with Oracle, Sun Microsystems, and AMD. Monica has worked with Java HotSpot VM optimizing the JIT compiler, the generated code, the JVM heuristics, and garbage collection and collectors. She is a regular speaker at various conferences and has several published articles on topics including garbage collection, the Java memory model, and others. Monica led Oracle's Garbage First Garbage Collector performance team, and was named a JavaOne Rock Star.

Poonam Parhar (Santa Clara, CA) is currently a JVM Sustaining Engineer at Oracle where her primary responsibility is to resolve customer-escalated problems against JRockit and HotSpot VMs. She loves debugging and troubleshooting problems and is always focused on improving the serviceability and supportability of the HotSpot VM. She has nailed down many complex garbage collection issues in the HotSpot VM and is passionate about improving the debugging tools and the serviceability of the product so as to make it easier to troubleshoot and fix garbage-collector-related issues in the HotSpotVM. She has made several contributions to the Serviceability Agent debugger and also developed a VisualVM plugin for it. She presented “VisualVM Plugin for the SA” at the JavaOne 2011 conference. In an attempt to help customers and the Java community, she shares her work experiences and knowledge through the blog she maintains at <https://blogs.oracle.com/poonam/>.

Bengt Rutisson (Stockholm, Sweden) is a JVM Engineer at Oracle, working on the HotSpot engineering team. He has worked on garbage collectors in JVMs for the past 10 years, first with the JRockit VM and the last six years with the HotSpot VM. Bengt is an active participant in the OpenJDK project, with many contributions of features, stability fixes, and performance enhancements.



Garbage First Overview

This chapter is an introduction to the Garbage First (or G1) garbage collector (GC) along with a historical perspective on the garbage collectors in the Java HotSpot Virtual Machine (VM), hereafter called just HotSpot, and the reasoning behind G1's inclusion in HotSpot. The reader is assumed to be familiar with basic garbage collection concepts such as young generation, old generation, and compaction. Chapter 3, "JVM Overview," of the book *Java™ Performance* [1] is a good source for learning more about these concepts.

Serial GC was the first garbage collector introduced in HotSpot in 1999 as part of Java Development Kit (JDK) 1.3.1. The Parallel and Concurrent Mark Sweep collectors were introduced in 2002 as part of JDK 1.4.2. These three collectors roughly correspond to the three most important GC use cases: "minimize memory footprint and concurrent overhead," "maximize application throughput," and "minimize GC-related pause times." One might ask, "Why do we need a new collector such as G1?" Before answering, let's clarify some terminology that is often used when comparing and contrasting garbage collectors. We'll then move on to a brief overview of the four HotSpot garbage collectors, including G1, and identify how G1 differs from the others.

Terminology

In this section, we define the terms *parallel*, *stop-the-world*, and *concurrent*. The term *parallel* means a multithreaded garbage collection operation. When a GC event activity is described as parallel, multiple threads are used to perform it. When a garbage

collector is described as parallel, it uses multiple threads to perform garbage collection. In the case of the HotSpot garbage collectors, almost all multithreaded GC operations are handled by internal Java VM (JVM) threads. One major exception to this is the G1 garbage collector, in which some background GC work can be taken on by the application threads. For more detail see Chapter 2, “Garbage First Garbage Collector in Depth,” and Chapter 3, “Garbage First Garbage Collector Performance Tuning.”

The term *stop-the-world* means that all Java application threads are stopped during a GC event. A stop-the-world garbage collector is one that stops all Java application threads when it performs a garbage collection. A GC phase or event may be described as stop-the-world, which means that during that particular GC phase or event all Java application threads are stopped.

The term *concurrent* means that garbage collection activity is occurring at the same time as the Java application is executing. A concurrent GC phase or event means that the GC phase or event executes at the same time as the application.

A garbage collector may be described by any one or a combination of these three terms. For example, a parallel concurrent collector is multithreaded (the parallel part) and also executes at the same time as the application (the concurrent part).

Parallel GC

Parallel GC is a parallel stop-the-world collector, which means that when a GC occurs, it stops all application threads and performs the GC work using multiple threads. The GC work can thus be done very efficiently without any interruptions. This is normally the best way to minimize the total time spent doing GC work relative to application work. However, individual pauses of the Java application induced by GC can be fairly long.

Both the young and old generation collections in Parallel GC are parallel and stop-the-world. Old generation collections also perform compaction. Compaction moves objects closer together to eliminate wasted space between them, leading to an optimal heap layout. However, compaction may take a considerable amount of time, which is generally a function of the size of the Java heap and the number and size of live objects in the old generation.

At the time when Parallel GC was introduced in HotSpot, only the young generation used a parallel stop-the-world collector. Old generation collections used a single-threaded stop-the-world collector. Back when Parallel GC was first introduced, the HotSpot command-line option that enabled Parallel GC in this configuration was `-XX:+UseParallelGC`.

At the time when Parallel GC was introduced, the most common use case for servers required throughput optimization, and hence Parallel GC became the default collector for the HotSpot Server VM. Additionally, the sizes of most Java heaps tended to be between 512MB and 2GB, which keeps Parallel GC pause times

relatively low, even for single-threaded stop-the-world collections. Also at the time, latency requirements tended to be more relaxed than they are today. It was common for Web applications to tolerate GC-induced latencies in excess of one second, and as much as three to five seconds.

As Java heap sizes and the number and size of live objects in old generation grew, the time to collect the old generation became longer and longer. At the same time, hardware advances made more hardware threads available. As a result, Parallel GC was enhanced by adding a multithreaded old generation collector to be used with a multithreaded young generation collector. This enhanced Parallel GC reduced the time required to collect and compact the heap.

The enhanced Parallel GC was delivered in a Java 6 update release. It was enabled by a new command-line option called `-XX:+UseParallelOldGC`. When `-XX:+UseParallelOldGC` is enabled, parallel young generation collection is also enabled. This is what we think of today as Parallel GC in HotSpot, a multithreaded stop-the-world young generation collector combined with a multithreaded stop-the-world old generation collector.

Tip

In Java 7 update release 4 (also referred to as Java 7u4, or JDK 7u4), `-XX:+UseParallelOldGC` was made the default GC and the normal mode of operation for Parallel GC. As of Java 7u4, specifying `-XX:+UseParallelGC` also enables `-XX:+UseParallelOldGC`, and likewise specifying `-XX:+UseParallelOldGC` also enables `-XX:+UseParallelGC`.

Parallel GC is a good choice in the following use cases:

1. Application throughput requirements are much more important than latency requirements.

A batch processing application is a good example since it is noninteractive. When you start a batch execution, you expect it to run to completion as fast as possible.

2. If worst-case application latency requirements can be met, Parallel GC will offer the best throughput. Worst-case latency requirements include both worst-case pause times, and also how frequently the pauses occur. For example, an application may have a latency requirement of “pauses that exceed 500ms shall not occur more than once every two hours, and all pauses shall not exceed three seconds.”

An interactive application with a sufficiently small live data size such that a Parallel GC's full GC event is able to meet or beat worst-case GC-induced latency requirements for the application is a good example that fits this use case. However, since the amount of live data tends to be highly correlated with the size of the Java heap, the types of applications falling into this category are limited.

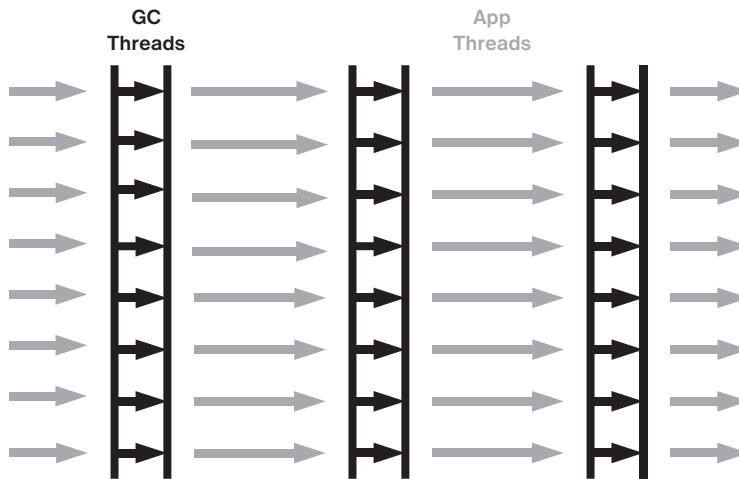


Figure 1.1 How Java application threads are interrupted by GC threads when Parallel GC is used

Parallel GC works well for applications that meet these requirements. For applications that do not meet these requirements, pause times can become excessively long, since a full GC must mark through the entire Java heap and also compact the old generation space. As a result, pause times tend to increase with increased Java heap sizes.

Figure 1.1 illustrates how the Java application threads (gray arrows) are stopped and the GC threads (black arrows) take over to do the garbage collection work. In this diagram there are eight parallel GC threads and eight Java application threads, although in most applications the number of application threads usually exceeds the number of GC threads, especially in cases where some application threads may be idle. When a GC occurs, all application threads are stopped, and multiple GC threads execute during GC.

Serial GC

Serial GC is very similar to Parallel GC except that it does all its work in a single thread. The single-threaded approach allows for a less complex GC implementation and requires very few external runtime data structures. The memory footprint is the lowest of all HotSpot collectors. The challenges with Serial GC are similar to those for Parallel GC. Pause times can be long, and they grow more or less linearly with the heap size and amount of live data. In addition, with Serial GC the long pauses are more pronounced, since the GC work is done in a single thread.

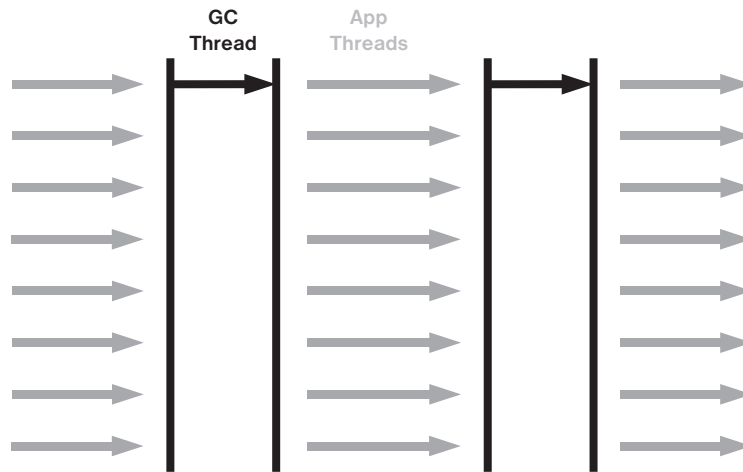


Figure 1.2 How Java application threads are interrupted by a single GC thread when Serial GC is used

Because of the low memory footprint, Serial GC is the default on the Java HotSpot Client VM. It also addresses the requirements for many embedded use cases. Serial GC can be explicitly specified as the GC to use with the `-XX:+UseSerialGC` HotSpot command-line option.

Figure 1.2 illustrates how Java application threads (gray arrows) are stopped and a single GC thread (black arrow) takes over to do the garbage collection work on a machine running eight Java application threads. Because it is single-threaded, Serial GC in most cases will take longer to execute a GC event than Parallel GC since Parallel GC can spread out the GC work to multiple threads.

Concurrent Mark Sweep (CMS) GC

CMS GC was developed in response to an increasing number of applications that demand a GC with lower worst-case pause times than Serial or Parallel GC and where it is acceptable to sacrifice some application throughput to eliminate or greatly reduce the number of lengthy GC pauses.

In CMS GC, young garbage collections are similar to those of Parallel GC. They are parallel stop-the-world, meaning all Java application threads are paused during young garbage collections and the garbage collection work is performed by multiple threads. Note that you can configure CMS GC with a single-threaded young generation collector, but this option has been deprecated in Java 8 and is removed in Java 9.

The major difference between Parallel GC and CMS GC is the old generation collection. For CMS GC, the old generation collections attempt to avoid long pauses in

application threads. To achieve this, the CMS old generation collector does most of its work concurrently with application thread execution, except for a few relatively short GC synchronization pauses. CMS is often referred to as mostly concurrent, since there are some phases of old generation collection that pause application threads. Examples are the initial-mark and remark phases. In CMS's initial implementation, both the initial-mark and remark phases were single-threaded, but they have since been enhanced to be multithreaded. The HotSpot command-line options to support multithreaded initial-mark and remark phases are `-XX:+CMSParallelInitialMarkEnabled` and `-XX:CMSParallelRemarkEnabled`. These are automatically enabled by default when CMS GC is enabled by the `-XX:+UseConcurrentMarkSweepGC` command-line option.

It is possible, and quite likely, for a young generation collection to occur while an old generation concurrent collection is taking place. When this happens, the old generation concurrent collection is interrupted by the young generation collection and immediately resumes upon the latter's completion. The default young generation collector for CMS GC is commonly referred to as ParNew.

Figure 1.3 shows how Java application threads (gray arrows) are stopped for the young GCs (black arrows) and for the CMS initial-mark and remark phases, and old generation GC stop-the-world phases (also black arrows). An old generation collection in CMS GC begins with a stop-the-world initial-mark phase. Once initial mark completes, the concurrent marking phase begins where the Java application threads are allowed to execute concurrently with the CMS marking threads. In Figure 1.3, the concurrent marking threads are the first two longer black arrows, one on top of the other below the "Marking/Pre-cleaning" label. Once concurrent marking completes, concurrent pre-cleaning is executed by the CMS threads, as shown by the two shorter black arrows under the "Marking/Pre-cleaning" label. Note that if there are enough available hardware threads, CMS thread execution overhead will not have much effect on the performance of Java application threads. If, however, the hardware threads are saturated or highly utilized, CMS threads will compete for CPU cycles with Java application threads. Once concurrent pre-cleaning completes, the stop-the-world remark phase begins. The remark phase marks objects that may have been missed after the initial mark and while concurrent marking and concurrent pre-cleaning execute. After the remark phase completes, concurrent sweeping begins, which frees all dead object space.

One of the challenges with CMS GC is tuning it such that the concurrent work can complete before the application runs out of available Java heap space. Hence, one tricky part about CMS is to find the right time to start the concurrent work. A common consequence of the concurrent approach is that CMS normally requires on the order of 10 to 20 percent more Java heap space than Parallel GC to handle the same application. That is part of the price paid for shorter GC pause times.

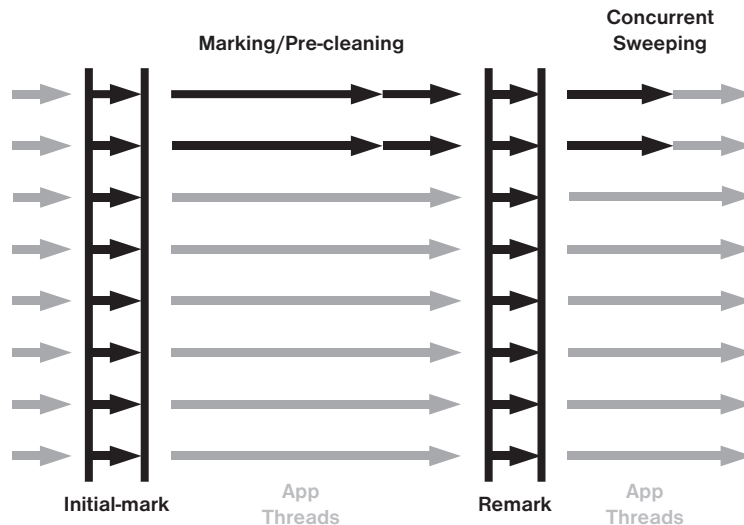


Figure 1.3 How Java application threads are impacted by the GC threads when CMS is used

Another challenge with CMS GC is how it deals with fragmentation in the old generation. Fragmentation occurs when the free space between objects in the old generation becomes so small or nonexistent that an object being promoted from the young generation cannot fit into an available hole. The CMS concurrent collection cycle does not perform compaction, not even incremental or partial compaction. A failure to find an available hole causes CMS to fall back to a full collection using Serial GC, typically resulting in a lengthy pause. Another unfortunate challenge associated with fragmentation in CMS is that it is unpredictable. Some application runs may never experience a full GC resulting from old generation fragmentation while others may experience it regularly.

Tuning CMS GC can help postpone fragmentation, as can application modifications such as avoiding large object allocations. Tuning can be a nontrivial task and requires much expertise. Making changes to the application to avoid fragmentation may also be challenging.

Summary of the Collectors

All of the collectors described thus far have some common issues. One is that the old generation collectors must scan the entire old generation for most of their operations such as marking, sweeping, and compacting. This means that the time to perform the work scales more or less linearly with the Java heap size. Another is that it must

be decided up front where the young and old generations should be placed in the virtual address space, since the young and old generations are separate consecutive chunks of memory.

Garbage First (G1) GC

The G1 garbage collector addresses many of the shortcomings of Parallel, Serial, and CMS GC by taking a somewhat different approach. G1 divides the heap into a set of regions. Most GC operations can then be performed a region at a time rather than on the entire Java heap or an entire generation.

In G1, the young generation is just a set of regions, which means that it is not required to be a consecutive chunk of memory. Similarly, the old generation is also just a set of regions. There is no need to decide at JVM launch time which regions should be part of the old or young generation. In fact, the normal operational state for G1 is that over time the virtual memory mapped to G1 regions moves back and forth between the generations. A G1 region may be designated as young and later, after a young generation collection, become available for use elsewhere, since young generation regions are completely evacuated to unused regions.

In the remainder of this chapter, the term *available region* is used to identify regions that are unused and available for use by G1. An available region can be used or designated as a young or old generation region. It is possible that after a young generation collection, a young generation region can at some future time be used as an old generation region. Likewise, after collection of an old generation region, it becomes an available region that can at some future time be used as a young generation region.

G1 young collections are parallel stop-the-world collections. As mentioned earlier, parallel stop-the-world collections pause all Java application threads while the garbage collector threads execute, and the GC work is spread across multiple threads. As with the other HotSpot garbage collectors, when a young generation collection occurs, the entire young generation is collected.

Old generation G1 collections are quite different from those of the other HotSpot collectors. G1 old generation collections do not require the entire old generation to be collected in order to free space in the old generation. Instead, only a subset of the old generation regions may be collected at any one time. In addition, this subset of old generation regions is collected in conjunction with a young collection.

Tip

The term to describe the collection of a subset of old generation regions in conjunction with a young collection is *mixed GC*. Hence, a mixed GC is a GC event in which all young generation regions are collected in addition to a subset of old generation regions. In other words, a mixed GC is a mix of young and old generation regions that are being collected.

Similar to CMS GC, there is a fail-safe to collect and compact the entire old generation in dire situations such as when old generation space is exhausted.

A G1 old generation collection, ignoring the fail-safe type of collection, is a set of phases, some of which are parallel stop-the-world and some of which are parallel concurrent. That is, some phases are multithreaded and stop all application threads, and others are multithreaded and execute at the same time as the application threads. Chapters 2 and 3 provide more detail on each of these phases.

G1 initiates an old generation collection when a Java heap occupancy threshold is exceeded. It is important to note that the heap occupancy threshold in G1 measures the old generation occupancy compared to the entire Java heap. Readers who are familiar with CMS GC remember that CMS initiates an old generation collection using an occupancy threshold applied against the old generation space only. In G1, once the heap occupancy threshold is reached or exceeded, a parallel stop-the-world initial-mark phase is scheduled to execute.

The initial-mark phase executes at the same time as the next young GC. Once the initial-mark phase completes, a concurrent multithreaded marking phase is initiated to mark all live objects in the old generation. When the concurrent marking phase is completed, a parallel stop-the-world remark phase is scheduled to mark any objects that may have been missed due to application threads executing concurrently with the marking phase. At the end of the remark phase, G1 has full marking information on the old generation regions. If there happen to be old generation regions that do not have any live objects in them, they can be reclaimed without any additional GC work during the next phase of the concurrent cycle, the cleanup phase.

Also at the end of the remark phase, G1 can identify an optimal set of old generations to collect.

Tip

The set of regions to collect during a garbage collection is referred to as a collection set (CSet).

The regions selected for inclusion in a CSet are based on how much space can be freed and the G1 pause time target. After the CSet has been identified, G1 schedules a GC to collect regions in the CSet during the next several young generation GCs. That is, over the next several young GCs, a portion of the old generation will be collected in addition to the young generation. This is the mixed GC type of garbage collection event mentioned earlier.

With G1, every region that is garbage collected, regardless of whether it is young or old generation, has its live objects evacuated to an available region. Once the live objects have been evacuated, the young and/or old regions that have been collected become available regions.

An attractive outcome of evacuating live objects from old generation regions into available regions is that the evacuated objects end up next to each other in the virtual

address space. There is no fragmented empty space between objects. Effectively, G1 does partial compaction of the old generation. Remember that CMS, Parallel, and Serial GC all require a full GC to compact the old generation and that compaction scans the entire old generation.

Since G1 performs GC operations on a per-region basis, it is suitable for large Java heaps. The amount of GC work can be limited to a small set of regions even though the Java heap size may be rather large.

The largest contributors to pause times in G1 are young and mixed collections, so one of the design goals of G1 is to allow the user to set a GC pause time goal. G1 attempts to meet the specified pause time goal through adaptive sizing of the Java heap. It will automatically adjust the size of the young generation and the total Java heap size based on the pause time goal. The lower the pause time goal, the smaller the young generation and the larger the total heap size, making the old generation relatively large.

A G1 design goal is to limit required tuning to setting a maximum Java heap size and specifying a GC pause time target. Otherwise, G1 is designed to dynamically tune itself using internal heuristics. At the time of writing, the heuristics within G1 are where most active HotSpot GC development is taking place. Also as of this writing, G1 may require additional tuning in some cases, but the prerequisites to building good heuristics are present and look promising. For advice on how to tune G1, see Chapter 3.

To summarize, G1 scales better than the other garbage collectors for large Java heaps by splitting the Java heap into regions. G1 deals with Java heap fragmentation with the help of partial compactions, and it does almost all its work in a multithreaded fashion.

As of this writing, G1 primarily targets the use case of large Java heaps with reasonably low pauses, and also those applications that are using CMS GC. There are plans to use G1 to also target the throughput use case, but for applications looking for high throughput that can tolerate longer GC pauses, Parallel GC is currently the better choice.

G1 Design

As mentioned earlier, G1 divides the Java heap into regions. The region size can vary depending on the size of the heap but must be a power of 2 and at least 1MB and at most 32MB. Possible region sizes are therefore 1, 2, 4, 8, 16, and 32MB. All regions are the same size, and their size does not change during execution of the JVM. The region size calculation is based on the average of the initial and maximum Java heap sizes such that there are about 2000 regions for that average heap size. As an example, for a 16GB Java heap with `-Xmx16g -Xms16g` command-line options, G1 will choose a region size of $16\text{GB}/2000 = 8\text{MB}$.

If the initial and maximum Java heap sizes are far apart or if the heap size is very large, it is possible to have many more than 2000 regions. Similarly, a small heap size may end up with many fewer than 2000 regions.

Each region has an associated remembered set (a collection of the locations that contain pointers into the region, shortened to RSet). The total RSet size is limited but noticeable, so the number of regions has a direct effect on HotSpot's memory footprint. The total size of the RSets heavily depends on application behavior. At the low end, RSet overhead is around 1 percent and at the high end 20 percent of the heap size.

A particular region is used for only one purpose at a time, but when the region is included in a collection, it will be completely evacuated and released as an available region.

There are several types of regions in G1. Available regions are currently unused. Eden regions constitute the young generation eden space, and survivor regions constitute the young generation survivor space. The set of all eden and survivor regions together is the young generation. The number of eden or survivor regions can change from one GC to the next, between young, mixed, or full GCs. Old generation regions comprise most of the old generation. Finally, humongous regions are considered to be part of the old generation and contain objects whose size is 50 percent or more of a region. Until a JDK 8u40 change, humongous regions were collected as part of the old generation, but in JDK 8u40 certain humongous regions are collected as part of a young collection. There is more detail on humongous regions later in this chapter.

The fact that a region can be used for any purpose means that there is no need to partition the heap into contiguous young and old generation segments. Instead, G1 heuristics estimate how many regions the young generation can consist of and still be collected within a given GC pause time target. As the application starts allocating objects, G1 chooses an available region, designates it as an eden region, and starts handing out memory chunks from it to Java threads. Once the region is full, another unused region is designated an eden region. The process continues until the maximum number of eden regions is reached, at which point a young GC is initiated.

During a young GC, all young regions, eden and survivor, are collected. All live objects in those regions are evacuated to either a new survivor region or to an old generation region. Available regions are tagged as survivor or old generation regions as needed when the current evacuation target region becomes full.

When the occupancy of the old generation space, after a GC, reaches or exceeds the initiating heap occupancy threshold, G1 initiates an old generation collection. The occupancy threshold is controlled by the command-line option `-XX:InitiatingHeapOccupancyPercent`, which defaults to 45 percent of the Java heap.

G1 can reclaim old generation regions early when the marking phase shows that they contain no live objects. Such regions are added to the available region set. Old regions containing live objects are scheduled to be included in a future mixed collection.

G1 uses multiple concurrent marking threads. In an attempt to avoid stealing too much CPU from application threads, marking threads do their work in bursts. They do as much work as they can fit into a given time slot and then pause for a while, allowing the Java threads to execute instead.

Humongous Objects

G1 deals specially with large object allocations, or what G1 calls “humongous objects.” As mentioned earlier, a humongous object is an object that is 50 percent or more of a region size. That size includes the Java object header. Object header sizes vary between 32- and 64-bit HotSpot VMs. The header size for a given object within a given HotSpot VM can be obtained using the Java Object Layout tool, also known as JOL. As of this writing, the Java Object Layout tool can be found on the Internet [2].

When a humongous object allocation occurs, G1 locates a set of consecutive available regions that together add up to enough memory to contain the humongous object. The first region is tagged as a “humongous start” region and the other regions are marked as “humongous continues” regions. If there are not enough consecutive available regions, G1 will do a full GC to compact the Java heap.

Humongous regions are considered part of the old generation, but they contain only one object. This property allows G1 to eagerly collect a humongous region when the concurrent marking phase detects that it is no longer live. When this happens, all the regions containing the humongous object can be reclaimed at once.

A potential challenge for G1 is that short-lived humongous objects may not be reclaimed until well past the point at which they become unreferenced. JDK 8u40 implemented a method to, in some cases, reclaim a humongous region during a young collection. Avoiding frequent humongous object allocations can be crucial to achieving application performance goals when using G1. The enhancements available in JDK 8u40 help but may not be a solution for all applications having many short-lived humongous objects.

Full Garbage Collections

Full GCs in G1 are implemented using the same algorithm as the Serial GC collector. When a full GC occurs, a full compaction of the entire Java heap is performed. This ensures that the maximum amount of free memory is available to the system. It is important to note that full GCs in G1 are single-threaded and as a result may introduce exceptionally long pause times. Also, G1 is designed such that full GCs are not expected to be necessary. G1 is expected to satisfy application performance

goals without requiring a full GC and can usually be tuned such that a full GC is not needed.

Concurrent Cycle

A G1 concurrent cycle includes the activity of several phases: initial marking, concurrent root region scanning, concurrent marking, remarking, and cleanup. The beginning of a concurrent cycle is the initial mark, and the ending phase is cleanup. All these phases are considered part of “marking the live object graph” with the exception of the cleanup phase.

The purpose of the initial-mark phase is to gather all GC roots. Roots are the starting points of the object graphs. To collect root references from application threads, the application threads must be stopped; thus the initial-mark phase is stop-the-world. In G1, the initial marking is done as part of a young GC pause since a young GC must gather all roots anyway.

The marking operation must also scan and follow all references from objects in the survivor regions. This is what the concurrent root region scanning phase does. During this phase all Java threads are allowed to execute, so no application pauses occur. The only limitation is that the scanning must be completed before the next GC is allowed to start. The reason for that is that a new GC will generate a new set of survivor objects that are different from the initial mark’s survivor objects.

Most marking work is done during the concurrent marking phase. Multiple threads cooperate to mark the live object graph. All Java threads are allowed to execute at the same time as the concurrent marking threads, so there is no pause in the application, though an application may experience some throughput reduction.

After concurrent marking is done, another stop-the-world phase is needed to finalize all marking work. This phase is called the “remark phase” and is usually a very short stop-the-world pause.

The final phase of concurrent marking is the cleanup phase. In this phase, regions that were found not to contain any live objects are reclaimed. These regions are not included in a young or mixed GC since they contain no live objects. They are added to the list of available regions.

The marking phases must be completed in order to find out what objects are live so as to make informed decisions about what regions to include in the mixed GCs. Since it is the mixed GCs that are the primary mechanism for freeing up memory in G1, it is important that the marking phase finishes before G1 runs out of available regions. If the marking phase does not finish prior to running out of available regions, G1 will fall back to a full GC to free up memory. This is reliable but slow. Ensuring that the marking phases complete in time to avoid a full GC may require tuning, which is covered in detail in Chapter 3.

Heap Sizing

The Java heap size in G1 is always a multiple of the region size. Except for that limitation, G1 can grow and shrink the heap size dynamically between `-Xms` and `-Xmx` just as the other HotSpot GCs do.

G1 may increase the Java heap size for several reasons:

1. An increase in size can occur based on heap size calculations during a full GC.
2. When a young or mixed GC occurs, G1 calculates the time spent to perform the GC compared to the time spent executing the Java application. If too much time is spent in GC according to the command-line setting `-XX:GCTimeRatio`, the Java heap size is increased. The idea behind growing the Java heap size in this situation is to allow GCs to happen less frequently so that the time spent in GC compared to the time spent executing the application is reduced.

The default value for `-XX:GCTimeRatio` in G1 is 9. All other HotSpot garbage collectors default to a value of 99. The larger the value for `GCTimeRatio`, the more aggressive the increase in Java heap size. The other HotSpot collectors are thus more aggressive in their decision to increase Java heap size and by default are targeted to spend less time in GC relative to the time spent executing the application.

3. If an object allocation fails, even after having done a GC, rather than immediately falling back to doing a full GC, G1 will attempt to increase the heap size to satisfy the object allocation.
4. If a humongous object allocation fails to find enough consecutive free regions to allocate the object, G1 will try to expand the Java heap to obtain more available regions rather than doing a full GC.
5. When a GC requests a new region into which to evacuate objects, G1 will prefer to increase the size of the Java heap to obtain a new region rather than failing the GC and falling back to a full GC in an attempt to find an available region.

References

- [1] Charlie Hunt and Binu John. *Java™ Performance*. Addison-Wesley, Upper Saddle River, NJ, 2012. ISBN 978-0-13-714252-1.
- [2] “Code Tools: jol.” OpenJDK, circa 2014. <http://openjdk.java.net/projects/code-tools/jol/>.



Index

A

Addresses, finding, 92
Age tables, 19
Aging, live objects, 18, 19
Allocating humongous objects, allocation path, 19, 21
Available regions, 8

B

Barriers, in RSets, 28–30, 35
Boolean command-line options, 145

C

Cards, definition, 25
Chunks. *See also* Regions
 cards, 25
 definition, 25–28
 global card table, 25–28
Class browser, 84–85
Class files, dumping, 84, 112
Class instances, displaying, 88
ClassDump, 106–108

Classes

 dumping, 106–108
 unloading, 153
Cleanup phase, G1 GC, 13
Cleanup stage of concurrent marking, 36
CLHSDB (Command-Line HotSpot Debugger). *See also* HSDB (HotSpot Debugger)
 command list, 102
 description, 100–102
 Java objects, examining, 101
 launching, 100
 VM data structures, examining, 101
CMS (Concurrent Mark Sweep) GC,
 pause times, 5–7
Coarse-grained bitmap, 25
Code root scanning, 43–44
Code Viewer, 95–97
Code Viewer panel, 120–122
Collection sets (CSets). *See* CSets (collection sets)
Command Line Flags, 98, 100

- Command-line options. *See also specific options*
 - boolean, 145
 - default values, 146
 - diagnostic options, enabling, 153–154
 - displaying, 98, 100
 - overtuning, causing evacuation failure, 59
 - types of, 145
- Compaction
 - CMS (Concurrent Mark Sweep) GC, 7
 - G1 GC, 10, 12
 - Parallel GC, 2–3
- Compute Reverse Pointers, 89–90
- Concurrent cycle, G1 GC, 13
- Concurrent garbage collection, definition, 2
- Concurrent Mark Sweep (CMS) GC, pause times, 5–7
- Concurrent marking
 - concurrent marking pre-write barriers, 35
 - description, 13
 - high remark times, 36
 - identifying allocated objects, 30
 - next bitmap, 30–33
 - NTAMS (next TAMS), 30–33, 37
 - overview, 33
 - previous bitmap, 30–33
 - pre-write barriers, 35
 - PTAMS (previous TAMS), 30–33, 37
 - SATB (snapshot-at-the-beginning), 30
 - SATB pre-write barriers, 35
 - TAMS (top-at-mark-starts), 30–33
- Concurrent marking, stages of
 - cleanup, 36
 - concurrent marking, 34–36
 - initial mark, 34
 - remark, 36
 - root region scanning, 34
- Concurrent marking cycle, triggering, 21
- Concurrent marking phase, tuning, 52–54
- Concurrent marking pre-write barriers, 35
- Concurrent marking stage of concurrent marking, 34–36
- Concurrent refinement threads, 42–43
 - logging, 148
 - maximum number, setting, 29
 - purpose of, 29
 - in RSets, 28–30
- Concurrent root region scanning, 13
- Constant pools, 70
- Copy to survivor, 18, 19
- Core files. *See also* Crash dump files
 - attaching to HSDB, 76
 - collecting, 108
 - Command Line Flags, 98, 100
 - command-line JVM options, displaying, 98, 100
 - creating, 108
 - description, 108
 - exploring with the SA Plugin, 118–119
 - Java threads, displaying, 119–120
 - JFR information, extracting, 107–108
 - JVM version, displaying, 98–99
 - multiple debugging sessions on, 113
 - opening with Serviceability Agent, 76–77
 - permanent generation statistics, printing, 103–104
 - process map, printing, 104
 - reading bits with the `/proc` interface, 111
 - remote debugging, 78–80, 114
 - thread stack trace, displaying, 119–120
- Core files, transported
 - debugging, 109–110
 - shared library problems, 109–110, 112
- Crash dump files. *See also* Core files
 - attaching HSDB to, 76
 - collecting, 108
 - creating, 108
 - description, 108
 - reading bits with the Windows Debugger Engine, 111
- CSets (collection sets)
 - definition, 9
 - minimum old CSet size, setting, 57–58
 - old regions per CSet, maximum threshold, 58
 - old regions per CSet, minimum threshold, 57–58

D

Deadlock Detection, 84, 86
Deadlocks
 detecting, 84, 86
 troubleshooting, 131–136
Debug messages, enable printing on
 Windows, 112
Debugging tools. *See* JDI (Java Debug Interface); Serviceability Agent
Deduplicating Java strings, 149–151
Dirty card queue, 29–30
DumpJFR, 107–108

E

Eden space, 11, 18–19
Efficiency, garbage collection. *See* GC efficiency
Ergonomics heuristics, dumping, 55
Evacuation, 47
Evacuation failures
 definition, 16
 duration of, determining, 37
 log messages, 53–54
 overview, 37–38
 to-space exhausted log message, 53–54
Evacuation failures, potential causes
 heap size, 59
 insufficient space in survivor regions, 60
 long-lived humongous objects, 59–60
 overtuning JVM command-line options, 59
External root region scanning, 42

F

Finalizable objects, printing details of, 103
FinalizerInfo, 103
Find Address in Heap, 92–93
Find Object by Query, 90–91
Find panel, 122–123
Find Pointer, 92, 122–123
Find Value in Code Cache, 92–93
Find Value in CodeCache (in VisualVM), 122–123

Find Value in Heap, 122–123
Fine RSets, 25
Fragmentation
 CMS (Concurrent Mark Sweep) GC, 7
 G1 GC, 10
 maximum allowable, 56–59
 per region, controlling, 58
Full GCs
 algorithm for, 12–13
 compacting the heap, 12–13
 evaluation failures, 37–38
 triggering, 37–38

G

G1 garbage collector. *See* G1 GC
G1 GC. *See also* Regions
 available regions, 8
 cleanup, 13
 collection cycles, 16–17
 compaction, 10, 12
 concurrent cycle, 13
 concurrent marking, 13
 concurrent root region scanning, 13
 description, 8–10
 enabling, 146
 fragmentation, 10
 free memory amount, setting, 146–147
 full GCs, 12–13
 heap sizing, 14
 initial-mark phase, 13
 marking the live object graph, 13
 pause times, 10
 remarking, 13
 threads, setting number of, 146
Garbage, maximum allowable, 56–59
Garbage collectors. common issues, 7–8
Garbage First garbage collector.
 See G1 GC
Garbage-collector-related issues, 92
GC efficiency
 definition, 17
 identifying garbage collection candidates, 22–24
GC Worker Other times too high, 48
GCLABs (GC local allocation buffers), 47
gcore tool, 108
Generations, definition, 24

H

Heap, occupancy threshold, adaptive adjustment, 151

Heap Parameters, 98

HeapDumper, 103

Heaps

- boundaries, displaying, 98–99, 101
- dividing into regions. *See* Regions
- dumping, 103
- free memory, setting, 152
- G1 concurrent cycle, initiating, 149
- liveness information, printing, 147
- raw memory contents, displaying, 94
- region size, setting, 146
- size, causing evacuation failure, 59
- size changes, logging, 152
- sizes, Parallel GC, 2–4
- sizing, G1 GC, 14

HotSpot Serviceability Agent.

- See* Serviceability Agent

Hprof files, creating, 149

HSDB (HotSpot Debugger). *See also*

- CLHSDB (Command-Line HotSpot Debugger)
- connecting to debug server, 78–80
- core files, attaching to, 76
- debugging modes, 74
- description, 72–80
- HotSpot processes, attaching to, 74–76
- launching, 72
- opening core files, 76–77
- remote debugging, 78–80

HSDB (HotSpot Debugger) tools and utilities

- addresses, finding, 92
- class browser, 84–85
- class instances, displaying, 88
- Code Viewer, 95–97
- Compute Reverse Pointers, 89–90
- Deadlock Detection, 84, 86
- displaying Java threads, 80–84
- dumping class files, 84, 112
- Find Address in Heap, 92–93
- Find Object by Query, 90–91
- Find Pointer, 92
- Find Value in Code Cache, 92–93
- garbage-collector-related issues, 92
- heap boundaries, displaying, 98–99

- Heap Parameters, 98
- Java Threads, 80–84
- JIT compiler-related problems, 92–93, 95–97
- liveness path, getting, 88–89, 91
- memory leaks, 87–89
- Memory Viewer, 94
- method bytecodes, displaying, 95–97, 120–122
- Monitor Cache Dump, 94
- Object Histogram, 87–89
- Object Inspector, 85–87
- out-of-memory problems, 87–89
- raw memory contents, displaying, 94
- reference paths, computing, 89–90
- synchronization-related issues, 95
- System Properties, 98–99
- VM Version Info, 98–99

Humongous continue regions, 12

Humongous objects

- allocation path, 19, 21
- causing evacuation failure, 59–60
- description, 12
- identifying, 19
- optimizing allocation and reclamation of, 20
- short-lived, optimizing collection of, 21

Humongous Reclaim, 50

Humongous regions

- contiguity, 21
- definition, 11
- description, 19–21
- examples of, 20–21

Humongous start regions, 12

I

IHOP. *See* `-XX:InitiatingHeapOccupancyPercent`

Initial-mark stage of concurrent marking, 34

Initial-mark phase, G1 GC, 13

`inspect` command, 101

J

Java application threads. *See* Threads

Java Debug Interface (JDI). *See* JDI (Java Debug Interface)

- Java heaps. *See* Heaps
 - Java objects. *See also* Humongous objects; Young collection pauses, live objects
 - examining, 101
 - ordinary object pointers, inspecting, 120–121
 - Java Stack Trace panel, 119–120
 - Java strings
 - deduplicating, 149–150
 - deduplication statistics, printing, 151
 - deduplication threshold, setting, 150
 - Java threads. *See* Threads
 - Java Threads panel, 119–120
 - `java.lang.OutOfMemoryError`, 123–130
 - `java.lang.OutOfMemoryError`, troubleshooting, 123–130
 - Java-level deadlocks, troubleshooting, 131–136
 - JavaScript Debugger (JSDB), 108
 - JavaScript interface to Serviceability Agent, 108
 - JDI (Java Debug Interface)
 - description, 113
 - SA Core Attaching Connector, 113
 - SA Debug Server Attaching Connector, 114
 - SAPID Attaching Connector, 113
 - JFR information, extracting, 107–108
 - JIT compiler-related problems, 92–93, 95–97
 - `jsadbugd` utility, 114
 - JSDB (JavaScript Debugger), 108
 - JVM version, displaying, 98–99
- L**
- Large object allocations. *See* Humongous objects
 - `LIBSAPROC_DEBUG` environment variable, 112
 - `libsaproc.so` binaries, 69
 - Licensed features, enabling, 154
 - Live object graph, marking, 13
 - Live objects, liveness factor per region, calculating, 22–24
 - Liveness information, printing, 147
 - Liveness path, getting, 88–89, 91
 - Load balancing, 47–48
 - Log messages, evacuation failures, 53–54
 - Logging
 - concurrent refinement threads, 148
 - heap size changes, 152
 - update log buffer, 29–30
- M**
- Marking the live object graph, 13
 - Marking threshold, setting, 54
 - Memory
 - free amount, setting, 146–147, 152
 - reserving for promotions, 147
 - Memory leaks, 87–89
 - Memory Viewer, 94
 - Method bytecodes, displaying, 95–97, 120–122
 - Mixed collection cycle
 - definition, 23
 - number of mixed collections, determining, 23–24
 - Mixed collection pause, 22–24
 - Mixed collection phase
 - ergonomics heuristics, dumping, 55
 - fragmentation, maximum
 - allowable, 56–59
 - fragmentation per region, controlling, 58
 - garbage, maximum allowable, 56–59
 - minimum old CSet size, setting, 57–58
 - old regions per CSet, maximum threshold, 58
 - old regions per CSet, minimum threshold, 57–58
 - reclaimable percentage threshold, 56–59
 - tuning, 54–56
 - Mixed GCs, 8, 147
 - Monitor Cache Dump, 94
 - Multithreaded parallel garbage collection, 1–2
 - Multithreaded reference processing, enabling, 62
 - Mutator threads, definition, 29
- N**
- Native methods *versus* `nmethods`, 44
 - Next bitmap, 30–33

Nmethods, 44

NTAMS (next TAMS), 30–33, 37

O

Object Histogram, 87–89, 105

Object histograms, collecting, 87–89, 105

Object Inspector, 85–87

Old generation, live objects

 age tables, 19

 aging, 19

Old-to-old references, 25

Old-to-young references, 25

Oop Inspector panel, 120–121

oops (ordinary object pointers), inspecting,
120–121

Out-of-memory problems, 87–89

P

Parallel GC

 compaction, 2–3

 definition, 1–2

 description, 2–4

 enabling, 2–3

 interrupting Java application threads, 4

 Java heap sizes, 2–4

 pause times, 4

 uses for, 3–4

Pauses

 Concurrent Mark Sweep (CMS) GC,
 5–7

 G1 GC, 10

 Parallel GC, 4

 Serial GC, 4

 time goal, setting, 151

 young collection, 18–19

.pdb files, setting the path to, 111

Performance tuning. *See* Tuning

Permanent generation statistics, printing,
103–104

PermStat, 103–104

Per-region-tables (PRTs), 25–26

PLABs, 18, 21, 152–153

PMap, 104

Postmortem analysis

 core files, 136–143

 Java HotSpot VM crashes, 136–143

Previous bitmap, 30–33

Previous TAMS (PTAMS), 30–33, 37

Pre-write barriers

 concurrent marking, 35

 in RSets, 35

Process maps, printing, 104

Processed buffers, 42–44

Processes

 debugging, 113

 exploring with the SA Plugin,
 118–119

 Java threads, displaying, 119–120

 JFR information, extracting, 107–108

 JVM version, displaying, 98–99

 permanent generation statistics,
 printing, 103–104

 process map, printing, 104

 reading bits with the /proc
 interface, 111

 reading bits with the Windows
 Debugger Engine, 111

 remote debugging, 114

 thread stack trace, displaying,
 119–120

Promotion failure. *See* Evacuation
failures

PRTs (per-region-tables), 25–26

PTAMS (previous TAMS), 30–33, 37

R

Raw memory contents, displaying, 94

Reading bits with the /proc interface, 111

Reclaimable percentage threshold, mixed
collection phase, 56–59

Reclaiming regions, 13

Reclamation, 47

Reference object types, 60

Reference paths, computing, 89–90

Reference processing

 enabling multithreaded reference
 processing, 62

 excessive weak references, correcting, 63

 observing, 60

 overhead, 36

 overview, 60

 reference object types, 60

 soft references, 63–65

- Region size
 - adaptive selection, overwriting, 18
 - average, 10
 - calculating, 10
 - determining, 18
 - setting, 146
 - Regions. *See also* Chunks; CSets (collection sets); RSets (remembered sets)
 - available, 8
 - candidates, identifying. *See* GC efficiency
 - concurrent root region scanning, 13
 - dividing heaps, 8–10
 - eden space, 11, 18–19
 - GC efficiency, 17
 - humongous, 11
 - maximum number of, 10–11
 - mixed GC, 8
 - mixing old generation with young collections, 8
 - reclaiming, 13
 - region size, 10
 - root, definition, 34
 - sorting. *See* GC efficiency
 - survivor fill capacity, 19
 - survivor space, 11, 18–19
 - types of, 11
 - Remark stage of concurrent marking, 36
 - Remark times, high, reasons for, 36
 - Remark phase, G1 GC, 13
 - Remembered sets (RSets). *See* RSets (remembered sets)
 - Remote debugging, 78–80, 114
 - Resizing, young generation, 18
 - RMI (Remote Method Invocation), 114
 - Root objects, definition, 34
 - Root references, collecting, 13
 - Root region scanning stage of concurrent marking, 34
 - RSets (remembered sets)
 - barriers, 28–30
 - coarse, 25
 - coarse-grained bitmap, 25
 - code root scanning, 43–44
 - concurrent marking pre-write barriers, 35
 - concurrent refinement threads, 28–30
 - definition, 11, 25
 - density, 25
 - fine, 25
 - number per region, 25
 - pre-write barriers, 35
 - printing a summary of, 147–148
 - processed buffers, 42–44
 - SATB pre-write barriers, 35
 - size, 11
 - sparse, 25
 - summarizing statistics, 44–47
 - visualizing potential improvements, 47
 - write barriers, 28–30
 - RSets (remembered sets), importance of
 - old-to-old references, 25
 - old-to-young references, 25
 - overview, 24–28
 - PRTs (per-region-tables), 25–26
- ## S
- SA_ALTRoot environment variable
 - description, 112
 - setting, 110
 - SA_IGNORE_THREADDB environment variable, 112
 - sa-jdi.jar binaries, 69, 113
 - SATB (snapshot-at-the-beginning), 30
 - SATB pre-write barriers, 35
 - sawindbg.dll binaries, 69
 - Scanning nmethods, 44
 - Serial GC. *See also* Parallel GC
 - definition, 4
 - description, 4–5
 - interrupting Java application threads, 5
 - pause times, 4
 - Serviceability Agent
 - binaries in the JDK, 69–70
 - components, 69
 - description, 68
 - enable printing of debug messages on Windows, 112
 - environment variables, 112
 - HotSpot data structures, 70–71
 - purpose of, 68
 - system properties, 111–112
 - version matching, 71

- Serviceability Agent, debugging tools.
 - See also* CLHSDB (Command-Line HotSpot Debugger); HSDB (HotSpot Debugger)
 - ClassDump, 106–108
 - DumpJFR, 107–108
 - extending, 115–117
 - finalizable objects, printing details of, 103
 - FinalizerInfo, 103
 - heap, dumping, 103
 - HeapDumper, 103
 - JavaScript interface to Serviceability Agent, 108
 - JFR information, extracting, 107–108
 - JSDB (JavaScript Debugger), 108
 - loaded classes, dumping, 106–108
 - Object Histogram, 105
 - object histograms, collecting, 105
 - permanent generation statistics, printing, 103–104
 - PermStat, 103–104
 - PMap, 104
 - process maps, printing, 104
 - SOQL (Structured Object Query Language), 106
 - SOQL queries, executing, 106
- Serviceability Agent, plugin for VisualVM
 - Code Viewer panel, 120–122
 - description, 117–118
 - Find panel, 122–123
 - Find Pointer utility, 122–123
 - Find Value in CodeCache utility, 122–123
 - Find Value in Heap utility, 122–123
 - installing, 118
 - Java Stack Trace panel, 119–120
 - Java Threads panel, 119–120
 - Oop Inspector panel, 120–121
 - ordinary object pointers, inspecting, 120–121
 - using, 118–119
 - utilities, 119–123
- Serviceability Agent, troubleshooting problems
 - Java-level deadlocks, 131–136
 - OutOfMemoryError, 123–130
 - postmortem analysis of core files, 136–143
 - postmortem analysis of HotSpot JVM crashes, 136–143
- Shared library problems with transported core files, 109–110
- Snapshot-at-the beginning (SATB), 30
- Soft references, 63–65
- SOQL (Structured Object Query Language), 106
- SOQL queries, executing, 106
- Sparse RSets, 25
- Stop-the-world garbage collection, 2–4
- sun.jvm.hotspot.debugger
 - .useProcDebugger, 111
- sun.jvm.hotspot.debugger
 - .useWindbgDebugger, 111
- sun.jvm.hotspot.debugger.windbg
 - .disableNativeLookup, 111
- sun.jvm.hotspot.debugger.windbg
 - .imagePath, 111
- sun.jvm.hotspot.debugger.windbg
 - .symbolPath, 111
- sun.jvm.hotspot.jdi
 - .SACoreAttachingConnector, 113
- sun.jvm.hotspot.jdi
 - .SAPIDAHachingConnector, 113
- sun.jvm.hotspot.loadLibrary
 - .DEBUG, 111
- sun.jvm.hotspot.runtime
 - .VM.disableVersionCheck, 111
- sun.jvm.hotspot.tools.jcore
 - .filter=<name of class> property, 106–108, 111
- sun.jvm.hotspot.tools.jcore
 - .outputDir=<output directory> property, 106–107
- sun.jvm.hotspot.tools.jcore
 - .PackageNameFilter
 - .pkgList=<list of packages> property, 106, 111
- sun.jvm.hotspot.tools.Tool, 115
- Survivor fill capacity, 19
- Survivor regions, insufficient space causing evacuation failure, 60
- Survivor space, 11, 18–19
- Synchronization-related issues, 95

T

TAMS (top-at-mark-starts), 30–33
 Tenuring threshold, live objects, 18, 19
 Termination, 47–48
 Thread count, increasing, 54
 Thread stack trace, displaying, 119–120
 Threads
 CMS (Concurrent Mark Sweep)
 GC, 6–7
 displaying, 80–84, 119–120
 setting number of, 146
 time ratio, setting, 148–149
 Threads, interrupting
 Parallel GC, 4
 Serial GC, 5
 Timed activities, variance in, 42
 TLABs, 17–18, 21, 153
 Top-at-mark-starts (TAMS), 30–33
 To-space exhausted log message,
 53–54. *See also* Evacuation failures
 To-space exhaustion. *See* Evacuation
 failures
 To-space overflow. *See* Evacuation failures
 Troubleshooting. *See* JDI (Java Debug
 Interface); Serviceability Agent
 Tuning
 concurrent marking phase, 52–54
 mixed collection phase, 54–56
 reclamation, 63–65
 young generations, 50–52

U

universe command, 101
 Update log buffer, 29–30
 userdump tool, 108

V

Version checking, disabling, 111
 Version matching, 71
 VisualVM. *See* Serviceability Agent
 VM data structures, examining, 101
 VM Version Info, 98–99
 VMStructs class, 70–71
 vmStructs.cpp file, 70–71

W

Weak references, excessive, correcting, 63
 Work stealing, 47–48
 Write barriers, RSets, 28–30

X

-XX:+ClassUnloadingWithConcurrent
 Mark, 36, 153
 -XX:+CMSParallelInitialMark
 Enabled, 6
 -XX:+CMSParallelRemarkEnabled, 6
 -XX:ConcGCThreads, 34–35, 54, 146
 -XX:G1ConcRefinementGreenZone,
 29, 148
 -XX:G1ConcRefinementRedZone,
 29, 148
 -XX:G1ConcRefinementThreads,
 29–30, 43
 -XX:G1ConcRefinementYellowZone,
 29, 148
 -XX:G1HeapRegionSize, 19, 59, 146
 -XX:G1HeapRegionSize=n, 18
 -XX:G1HeapWastePercent,
 23, 146–147
 -XX:G1MaxNewSizePercent,
 17, 50–52
 -XX:G1MixedGCCCountTarget,
 23, 57, 147
 -XX:G1MixedGCLiveThreshold
 Percent, 58
 -XX:G1NewSizePercent, 17, 50–52
 -XX:+G1PrintRegionLiveness
 Info, 147
 -XX:G1ReservePercent, 60, 147
 -XX:+G1SummarizeRSetStats,
 44–47, 147
 -XX:G1SummarizeRSetStats
 Period, 148
 -XX:+G1TraceConcRefinement, 148
 -XX:+G1UseAdaptiveConc
 Refinement, 148
 -XX:+G1UseAdaptiveIHOP, 151
 -XX:GCTimeRatio, 148–149
 -XX:+HeapDumpAfterFullGC, 149
 -XX:+HeapDumpBeforeFullGC, 149

- XX:InitiatingHeapOccupancyPercent
 - default value, 11
 - description, 149
 - heap occupancy percentage, 22
 - occupancy threshold, default, 22
 - occupancy threshold, setting, 30
 - overview, 22–24
- XX:InitiatingHeapOccupancyPercent=n, 52–54
- XX:MaxGCPauseMillis, 17, 50–52, 151
- XX:MaxHeapFreeRatio, 152
- XX:MaxTenuringThreshold, 19
- XX:MinHeapFreeRatio, 152
- XX:ParallelGCThreads, 29–30, 35, 43, 54
- XX:+ParallelRefProcEnabled, 62–63
- XX:+PrintAdaptiveSizePolicy, 55, 59, 152
- XX:+PrintGCDetails, 18, 27, 60–61
- XX:PrintGCTimeStamps, 27
- XX:+PrintReferenceGC, 61–63
- XX:+PrintStringDeduplicationStatistics, 151
- XX:+ResizePLAB, 152–153
- XX:+ResizeTLAB, 153
- XX:SoftRefLRUPolicyMSPerMB=1000, 63–64
- XX:StringDeduplicationAgeThreshold, 150
- XX:TargetSurvivorRatio, 19
- XX:+UnlockCommercialFeatures, 154
- XX:+UnlockDiagnosticVMOptions, 147, 153–154
- XX:+UnlockExperimentalVMOptions, 154
- XX:+UseConcurrentMarkSweepGC, 6
- XX:+UseG1GC, 27
- XX:+UseG1GC, 146
- XX:+UseParallelGC, 2–3
- XX:+UseParallelOldGC, 2–3
- XX:+UseSerialGC, 5
- XX:+UseStringDeduplication, 149–150

Y

- Young collection pauses
 - description, 18–19
 - eden regions, 18–19
 - survivor regions, 18–19
 - triggering, 18
- Young collection pauses, live objects
 - aging, 18, 19
 - copy to survivor, 18, 19
 - identifying, 22–24
 - liveness factor per region, calculating, 22–24
 - survivor fill capacity, 19
 - tenuring, 18, 19
 - tenuring threshold, 18, 19
- Young collections
 - concurrent marking phase, tuning, 52–54
 - evacuation failure, log messages, 53–54
 - increasing thread count, 54
 - marking threshold, setting, 54
 - phases, 39. *See also specific phases*
- Young collections, parallel phase
 - activities outside of GC, 48
 - code root scanning, 43–44
 - code sample, 39–40
 - concurrent refinement threads, 42–43
 - definition, 39
 - evacuation, 47
 - external root region scanning, 42
 - load balancing, 47–48
 - processed buffers, 42–44
 - reclamation, 47
 - RSets, and processed buffers, 42–44
 - RSets, summarizing statistics, 44–47
 - scanning nmethods, 44
 - start of parallel activities, 41
 - summarizing parallel activities, 48
 - termination, 47–48

- variance in timed activities, 42
- work stealing, 47–48
- Young collections, serial phase
 - activities, 48–50
 - definition, 39
- Young generation
 - description, 17–18
 - initial generation size, setting, 50–52
 - maximum growth limit, setting, 50–52
 - pause time goal, setting, 50–52
 - resizing, 18
 - size, calculating, 17–18
 - tuning, 50–52
- Yuasa, Taiichi, 30