



# Virtualizing and Tuning Large-Scale Java Platforms

Emad Benjamin



# **Virtualizing and Tuning Large-Scale Java Platforms**

VMware Press is the official publisher of VMware books and training materials, which provide guidance on the critical topics facing today's technology professionals and students. Enterprises, as well as small- and medium-sized organizations, adopt virtualization as a more agile way of scaling IT to meet business needs. VMware Press provides proven, technically accurate information that will help them meet their goals for customizing, building, and maintaining their virtual environment.

With books, certification and study guides, video training, and learning tools produced by world-class architects and IT experts, VMware Press helps IT professionals master a diverse range of topics on virtualization and cloud computing. It is the official source of reference materials for preparing for the VMware Certified Professional Examination.

VMware Press is also pleased to have localization partners that can publish its products into more than 42 languages, including Chinese (Simplified), Chinese (Traditional), French, German, Greek, Hindi, Japanese, Korean, Polish, Russian, and Spanish.

For more information about VMware Press, please visit **[vmwarepress.com](http://vmwarepress.com)**.

# vmware® PRESS



[pearsonitcertification.com/vmwarepress](http://pearsonitcertification.com/vmwarepress)

Complete list of products • Podcasts • Articles • Newsletters

**VMware® Press** is a publishing alliance between Pearson and VMware, and is the official publisher of VMware books and training materials that provide guidance for the critical topics facing today's technology professionals and students.

With books, certification and study guides, video training, and learning tools produced by world-class architects and IT experts, VMware Press helps IT professionals master a diverse range of topics on virtualization and cloud computing, and is the official source of reference materials for completing the VMware certification exams.



Make sure to connect with us!  
[informit.com/socialconnect](http://informit.com/socialconnect)

vmware®

PEARSON  
IT CERTIFICATION

Safari®  
Books Online

*This page intentionally left blank*

# **Virtualizing and Tuning Large-Scale Java Platforms**

Emad Benjamin

**vmware® PRESS**

Upper Saddle River, NJ • Boston • Indianapolis • San Francisco  
New York • Toronto • Montreal • London • Munich • Paris • Madrid  
Capetown • Sydney • Tokyo • Singapore • Mexico City

## Virtualizing and Tuning Large-Scale Java Platforms

Copyright © 2014 VMware, Inc.

Published by Pearson plc

Publishing as VMware Press

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.

Library of Congress Control Number: 2013920560

ISBN-13: 978-0-13-349120-3

ISBN-10: 0-13-349120-X

Text printed in the United States on recycled paper at RR Donnelly in Crawfordsville, Indiana.

First Printing, December 2013

All terms mentioned in this book that are known to be trademarks or service marks have been appropriately capitalized. The publisher cannot attest to the accuracy of this information. Use of a term in this book should not be regarded as affecting the validity of any trademark or service mark.

VMware terms are trademarks or registered trademarks of VMware in the United States, other countries, or both.

### Warning and Disclaimer

Every effort has been made to make this book as complete and as accurate as possible, but no warranty or fitness is implied. The information provided is on an “as is” basis. The authors, VMware Press, VMware, and the publisher shall have neither liability nor responsibility to any person or entity with respect to any loss or damages arising from the information contained in this book.

The opinions expressed in this book belong to the author and are not necessarily those of VMware.

### Corporate and Government Sales

VMware Press offers excellent discounts on this book when ordered in quantity for bulk purchases or special sales, which may include electronic versions and/or custom covers and content particular to your business, training goals, marketing focus, and branding interests. For more information, please contact U.S. Corporate and Government Sales, (800) 382-3419, [corpsales@pearsontechgroup.com](mailto:corpsales@pearsontechgroup.com).

For sales outside the United States, please contact International Sales, [international@pearsoned.com](mailto:international@pearsoned.com).

### ASSOCIATE PUBLISHER

David Dusthimer

### ACQUISITIONS EDITOR

Joan Murray

### VMWARE PRESS

### PROGRAM MANAGER

Anand Sundaram

### SENIOR DEVELOPMENT EDITOR

Christopher Cleveland

### MANAGING EDITOR

Sandra Schroeder

### SENIOR PROJECT EDITOR

Tonya Simpson

### COPY EDITOR

Keith Cline

### PROOFREADER

Debbie Williams

### BOOK DESIGNER

Gary Adair

### COVER DESIGNER

Chuti Prasertsith

### COMPOSITOR

Bumpy Design

### EDITORIAL ASSISTANT

Vanessa Evans

*I dedicate this book to my beloved wife, Christine,  
and our two beautiful boys, Anthony and Adrian.  
They fill my life with plenty of joy and inspiration.  
Thank you Lord for all your blessings.*



*This page intentionally left blank*

# Contents

## Preface xv

## Chapter 1 Introduction to Large-Scale Java Platforms 1

Large-Scale Java Platform Categories	1
Large-Scale Java Platform Trends and Requirements	2
Compute-Resource Consolidation	2
JVM Instance Consolidation	3
Elasticity and Flexibility	3
Performance	4
Large-Scale Java Platform Technical Considerations	4
Theoretical and Practical Limits of Java Platforms	4
NUMA	9
Most Common JVM Size Found in Production Environments	15
Horizontal Scaling Versus Vertical Scaling of JVMs and VMs	16
Summary	20

## Chapter 2 Modern Scalable Data Platforms 21

SQLFire Topologies	24
Client/Server Topology	24
Peer-to-Peer Topology	27
Redundancy Zones	28
Global Multisite Topology	28
SQLFire Features	30
Server Groups	32
Partitioning	34
Redundancy	37
Colocation	38
Disk Persistence	39
Transactions	41
Cache Plug-In	46
Listeners	47
Writers	50
Asynchronous Listeners	52
DBSynchronizer	54
SQLF Commands and DDLUtils	57
Active-Active Architectures and Modern Data Platforms	57
Chapter Summary	61

## **Chapter 3 Tuning Large-Scale Java Platforms 63**

- GC Tuning Approach 70
  - Step A: Young Generation Tuning 71
  - Step B: Old Generation Tuning 76
  - Step C: Survivor Spaces Tuning 78
- Chapter Summary 78

## **Chapter 4 Designing and Sizing Large-Scale Java Platforms 79**

- Designing and Sizing a New Environment for a Virtualized Large-Scale Java Platform 79
  - Step1: Establishing Your Current Production Load Profile 80
  - Step 2: Establish a Benchmark 82
  - Step 3: Size the Production Environment 95
- Sizing vFabric SQLFire Java Platforms: Category 2 Workloads 96
  - Step A: Determine Entity Groups 97
  - Step B: Determine the Memory Size of the Data Fabric 100
  - Step C: Establish Building Block VM and JVM Size and How Many vFabric SQLFire Members Are Needed 105
  - Understanding the Internal Memory Sections of HotSpot JVM 106
  - Understanding NUMA Implications on Sizing Large VMs and JVMs 108
  - vFabric SQLFire Sizing Example 112
- Chapter Summary 119

## **Chapter 5 Performance Studies 121**

- SQLFire Versus RDBMS Performance Study 121
  - Performance Results 123
  - Summary of Findings 126
- The Olio Workload on tc Server and vSphere Performance Study 127
  - Looking at the Results 127
- SpringTrader Performance Study 131
  - Application and Data Tier vSphere Configurations 133
  - The SpringTrader Performance Study Results 137
- Performance Differences Between ESXi 3, 4.1, and 5 139
  - CPU Scheduling Enhancements 140
  - Memory Enhancements 140
- vSphere 5 Performance Enhancements 142
- Chapter Summary 143

---

**Chapter 6 Best Practices 145**

Enterprise Java Applications on vSphere Best Practices (Category 1)	148
VM Sizing and Configuration Best Practices	148
vCPU for VM Best Practices	149
VM Memory Size Best Practices	150
VM Timekeeping Best Practices	156
Vertical Scalability Best Practices	156
Horizontal Scalability, Clusters, and Pools Best Practices	158
Inter-Tier Configuration Best Practices	160
High-Level vSphere Best Practices	165
SQLFire Best Practices and SQLFire on vSphere Best Practices (Category 2 JVM Workload Best Practices)	166
SQLFire Best Practices	168
vFabric SQLFire Best Practices on vSphere	173
Category 3 Workloads Best Practices	181
IBM JVM and Oracle jRockit JVMs	181
GC Policy Selection	184
IBM GC Choices	186
Oracle jRockit GC Policies	187
Chapter Summary	187

**Chapter 7 Monitoring and Troubleshooting Primer 189**

Open a Support-Request Ticket	191
Collecting Metrics from vCenter	191
Troubleshooting Techniques for vSphere with esxtop	195
Java Troubleshooting Primer	198
Troubleshooting Java Memory Problems	202
Troubleshooting Java Thread Contentions	203
Chapter Summary	204

**Appendix FAQs 205****Glossary 229****Index 233**

## Best Practices

Best Practice 1: Common Distributed Data Platform	24
Best Practice 2: Client/Server Topology	26
Best Practice 3: Peer-to-Peer Multihomed Machines	27
Best Practice 4: Multisite	29
Best Practice 5: Use Server Groups	33
Best Practice 6: Horizontal Partitioning	37
Best Practice 7: Redundancy	38
Best Practice 8: Colocation	38
Best Practice 9: Disk Persistence	40
Best Practice 10: Transactions	45
Best Practice 11: RowLoader	47
Best Practice 12: Listeners	49
Best Practice 13: Writers	51
Best Practice 14: Asynchronous Listeners	53
Best Practice 15: DBSynchronizer	55
Best Practice 16: VM Sizing and VM-to-JVM Ratio Through a Performance Load Test	149
Best Practice 17: VM vCPU CPU Overcommit	149
Best Practice 18: VM vCPU, Do Not Oversubscribe to CPU Cycles That You Don't Really Need	150
Best Practice 19: VM Memory Sizing	152
Best Practice 20: Set Memory Reservation for VM Memory Needs	154
Best Practice 21: Use of Large Pages	154
Best Practice 22: Use an NTP Source	156
Best Practice 23: Hot Add or Remove CPU/Memory	157
Best Practice 24: Use vSphere Host Clusters	158
Best Practice 25: Use Resource Pools	159
Best Practice 26: Use Affinity Rules	159
Best Practice 27: Use vSphere-Aware Load Balancers	160
Best Practice 28: Establish Appropriate Thread Ratios That Prevents Bottlenecks (HTTP threads:Java threads:DB Connections Ratio)	160
Best Practice 29: Apache Web Server Sizing	161
Best Practice 30: Load-Balancer Algorithm Choice and VM Symmetry	164
Best Practice 31: vSphere 5.1	165
Best Practice 32: vSphere Networking	165
Best Practice 33: vSphere Storage	166
Best Practice 34: vSphere Host	166
Best Practice 35: JVM Version	168
Best Practice 36: Use Parallel and CMS GC Policy Combination	168
Best Practice 37: Set Initial Heap Equal to Maximum Heap	170
Best Practice 38: Disable Calls to <code>System.gc()</code>	171

---

Best Practice 39: New Generation Size	171
Best Practice 40: Using 32-Bit Addressing in a 64-Bit JVM	171
Best Practice 41: Stack Size	172
Best Practice 42: Perm Size	172
Best Practice 43: Table Placements in a JVM	172
Best Practice 44: Enable Hyperthreading and Do Not Overcommit CPU	173
Best Practice 45: CPU Cache Sharing	175
Best Practice 46: vFabric SQLFire Member Server, JVM and VM Ratio	175
Best Practice 47: VM Placement	175
Best Practice 48: Set VM Memory Reservation	175
Best Practice 49: vMotion, DRS Cluster, and vFabric SQLFire Server	176
Best Practice 50: VMware HA and vFabric SQLFire	177
Best Practice 51: Guest OS	177
Best Practice 52: Physical NIC	177
Best Practice 53: Virtual NIC	178
Best Practice 54: Troubleshooting SYN Cookies	179
Best Practice 55: Storage	181

*This page intentionally left blank*

# Preface

This book is the culmination of 9 years of experience in running Java on VMware vSphere, both at VMware and at many of VMware's customers. In fact, many VMware customers run enterprise-critical Java applications on VMware vSphere and have achieved better total cost of ownership (TCO) and service level agreements (SLAs). In my first book, *Enterprise Java Applications Architecture on VMware*, the topic of Java virtualization was covered well, both from a high-level architecture perspective and with in-depth technical chapters on sizing and best practices. To keep that first book more affordable, a decision was made to hold back some of the chapters for a second book, what you are reading now. These two books are complementary in many ways. The first book has a few high-level chapters for architects, engineers, and managers considering Java virtualization for the first time and asking the high-level question "why." This book is all about how and what to tune for optimal performance.

Limiting the scope of the first book was a good idea; the book was thus made available quickly for those launching their first Java virtualization projects. It has been almost 2 years since the release date of that first book, and since then nearly 300 customer interactions have helped to further analyze the guidance offered. Some of these interactions have included large-scale Java platforms of significant scale and have significantly contributed to the greater level of detail in this book. This book discusses in detail the sizing and tuning of both small-scale and large-scale virtualized Java platforms—100 Java Virtual Machines (JVMs) to 10,000 JVMs and a JVM heap range of 1 to 128GB. This recent experience, combined with my 15 years of tuning Java platforms, is presented in this book in such way to summarize what is most practical and immediately applicable to the vast majority of Java workload types. You can retrofit the advice, deployment configuration, and garbage collection (GC) tuning knowledge gained from this book to effectively combat GC poor behavior or to design and size your Java platform overall. The best practices highlighted throughout this book apply to physical environments, virtual environments, or both.

## Motivation for Writing This Book

I have spent the past 9 years at VMware in various capacities ensuring that all internal enterprise Java applications were virtualized to showcase to VMware customers the benefits of the approach. In that time, I came to believe that a lot of the best practices that we learned from empirical evidence in production environments should be shared with the VMware community. I received lots of feedback requesting that I document many of the lessons learned and the various tips and tricks needed to successfully run enterprise Java applications on VMware. This served as the motivation for the first book, *Enterprise Java Applications Architecture on VMware* (<https://www.createspace.com/3632131>).



Continuing on from the motivation of the first book, this book (the second book) focuses on what to tune, how far you can tune it, and how large virtualized Java platforms can be. In essence, the first book had a reasonable mix of the “why virtualize” and “what/how to virtualize.” In contrast, this book examines “how large of a scale you can virtualize and how far you can drive the platform tuning.”

It was quite exciting to write the first book, as we were trying to let the broader VMware customer base know that Java virtualization absolutely works and provides significant advantages. In this current book, we want to help those customers who are saying, “Now help me take it to the next level of scale.” We have spent the past 2 years helping customers virtualize many large-scale JVM platforms, some as big as 10,000 JVMs, and others in the big data platform space (with multiple terabytes of data kept in memory within a set of clustered JVMs). Before you dive into this book, though, remember this: Although this book presents many best practices, these practices represent optimal configuration guidance; they are not mandatory requirements. In our experience, we have found that most enterprise Java applications virtualize readily without having to worry about too many specific configurations. In fact, of any enterprise-level production application, Java applications are prime “low-hanging fruit” candidates for virtualization. By sharing the lessons we learned, we hope that you can avoid some of the pitfalls we encountered in our efforts to virtualize large-scale Java platforms.

Wanting to cover how best to deploy Java platforms in virtual environments (while also addressing misconceptions that the virtual platform is the problem), we built best practices that apply equally to both physical and virtual environments. By design, this book contains sections that cover best practices for physical and virtualized Java platforms, so as to enable customers to correct any problems on their physical Java platform before they virtualize. Of course, this is not mandatory; customers may choose to keep the legacy aspect of their physical Java deployment as they migrate to virtualization. However, at least they have been made aware of the design and deployment deficiencies of their physical Java platform should they wish to correct it in the future.

This is an important exercise to go through, allowing us to highlight the problem was actually customers’ own physical environments. Customers could thus understand the cost of maintaining the legacy aspect of their physical Java platform. For example, we often discover that many Java physical platforms were poorly architected with the wrong deployment topology, many times with sprawling thousands of unnecessary JVMs.

When we speak with customers, we walk them through the best practices and ensure that these environments are sized and tuned correctly, regardless of whether they leave their Java applications on physical environments or migrate to virtual. Again, customers can choose to ignore our prescriptions and deploy the legacy aspect of their Java physical

platform onto the virtual equivalent without much change or intrusion on the codebase and platform. However, customers these days are pretty cognizant of the value of the best practices we (and many others) have embraced to improve the Java deployment paradigm while migrating to virtual platforms.

The lessons learned fall into the following general categories:

- Things will go wrong in production; it is just a matter of when. So, you want to meticulously consider what could go wrong and have a roll-forward and a rollback plan. The planning exercise helps to further solidify the QA test plan. Note that this is not specific to a virtualized environment. In fact, it is an equally stringent requirement whether you are dealing with a physical or a virtual infrastructure. However, the reality is that virtualization gives you the mechanisms to quickly deal with issues (in contrast to a physical case in which you are restricted to the amount of flexibility you have to move around your compute resources).
- Enterprise Java applications are the low-hanging fruit when it comes to virtualization.
- Everyone operated in various silos at each of the Java tiers and did not necessarily speak the same language in terms of technology and organizational logistics. This was certainly the *modus operandi* under the old physical (nonvirtualized) paradigm, with these technology and organizational silos having been formed over the past decades. However, cross-team collaboration was a big part of virtualizing Java on VMware; it drove a lot of the teams to talk to each other to facilitate a best-of-breed design. Teams from both application development and operations came to the table many times.
- Customers sometimes seek to rationalize the legacy aspects of their environment. As a consequence, customers pay additional administrative cost associated with sprawling JVMs in the physical environment; if not remedied, these costs carry over into the virtualized system. For example, do you really need those 1GB heap space 5000 JVMs? Couldn't they be consolidated? Absolutely, they can be, and we show you how you can save by reducing your licensing costs and improving administration (because you will have fewer JVMs to manage).
- Performance issues. Customers often race to the conclusion that any problem must be a virtualization issue or a GC issue. In reality, though, virtualization is not the issue, but the GC may sometimes be an issue. If a GC issue exists, though, it is not specific to virtualization, and in fact the issue is almost always equally present in the physical deployment.

- Big in-memory databases on physical Java platforms (1TB memory cluster, really!)? Absolutely, if the prime objective is to service transactions at any cost, but at the highest speed possible, this is the right architecture for you. I found that many customers were skeptical about this. The fact that they attempted to size these types of environments without regard for the underlying platform gave them a poor start. You must pay attention to the server machine architecture when sizing these data platforms (as discussed later). The other poor practice I found with some customers is that they attempted to size these environments to have 30 or so JVMs. Well, that is not the right approach, because maintaining chatter between that many JVMs at a high rate can make latency worse overall. Quite simply, these are latency-sensitive memory-bound workloads and perform better using a deployment paradigm of fewer larger JVMs. If you were to compare the performance of an in-memory database system that had 30 JVMs versus one with 8 much larger JVMs, the configuration with 8 larger JVMs would be better. Of course, the caveat here is that the larger JVMs are sized correctly for NUMA optimization and the appropriate GC tuning has been applied.
- Can big in-memory databases really perform when virtualized? Over the past couple of years, I have seen an increase in customers virtualizing Java application servers. Specifically, large-scale platforms with thousands of JVMs are being virtualized. One unique category of workload is in-memory data management systems that require terabytes of memory and are latency sensitive. With these in-memory data clusters, we find that although there are fewer JVMs, they do tend to be of large heap space, ranging in JVM size from 8 to 128GB (and usually fewer than 12 JVMs). Of course, there is nothing magical about the number 12, it could be as low as 3, or as high as 30. However, the more JVMs you have, the more potential latency issues you risk because of the additional network hops. Later in the book, you will learn how to size and tune these workloads.

Many Java application developers know the development process well, know how to write Java code, and know how to tune the JVM. However, too often that information stays with developers and is not shared with (or translated for) application administrators. Often, the skills needed to run Java platforms are split between Java developers and administrators, without a single person understanding both purviews. This silo-ing of understanding is changing, though, as more individuals begin to understand how to write Java code, deploy it, tune the JVM, and recognize the full breadth of virtualization and the intricacies of the server hardware architecture. So, another goal of this book is to encourage readers to follow this career path as they develop this skill-set profile.

I sincerely hope this book helps those with backgrounds in Java development, operations infrastructure, and virtualization. You can use this book both as a guide to combat day-to-day situations and as an aid to help you compose a strategic architectural map of your Java platform.

## **Prerequisites**

This book assumes a high-level understanding of Java, JVM GC, server hardware architecture, and virtualization technologies. The information herein relates to running large-scale Java platforms. Although you might want to brush up on virtualization before delving into the material in this book, most senior Java specialists will quickly learn enough about virtualization and virtualizing Java applications from this book alone. In fact, by learning the answer to the following question, you will gain enough background on virtualization to continue through this book. This is the question asked on day one by folks new to virtualizing Java: “Is Java both operating system and hypervisor independent?”

This text assumes that the reader has some background in the Java language and specifically JVM architecture. The book does its best to summarize the JVM architecture, and specific JVM tunings, but it is not a replacement for a book dedicated to Java tuning. Suffice it to say, vSphere administrators who are new to JVM tuning should be able to learn enough from this book to hold technical conversations with their Java counterparts. In addition, vSphere and Java administrators can apply the tuning advice in this book and modify it as it applies to their environment. The various chapters on JVM tuning advice have been written in such a way that they are less overwhelming than those found in other Java books, enabling vSphere and the Java administrators to have a quick and effective go-forward design and tuning strategy. Many VMware customers running Java have applied the tuning parameters discussed in this book and have gained immediate performance improvements.

## **What You Need to Know First**

The fact that you are reading this book means that you are halfway to correcting your Java platform. Perhaps you have already concluded that tuning Java platforms cannot be ignored/minimized. However, even if you have just dabbled with VMware virtualization and Java tuning to some degree on physical systems, you are ready for this book. As a refresher (or an introduction for those new to the material), the following sections briefly introduce important concepts related to virtualizing and tuning large-scale Java platforms.

### **Is a 4GB Java Heap the New 1GB? Why?**

In the past 2 years, I have conducted more than 290 customer calls and workshops where it was evident that 40% of workloads running on Java platforms were deployed on JVMs that were 1 to 4GB in size. I continue to see a huge number of less than 1GB JVMs, approximately another 40% within the customer base I interact with. The remaining 20% varies from 4 to 360GB. Yes, a 360GB JVM is due to a monitoring system that cannot horizontally scale out and so the customer is forced to have a single JVM. Although this might seem unbelievable, it is the reality of Java production platforms today. However, the JVMs with 1GB of heap cause a sprawl of JVMs instances, and that becomes its own management headache. For example, you might want to service 1TB of total heap space, which would then mean thousands of JVM instances if you allow only a 1GB JVM heap. How can that possibly make any sense? Couldn't you get the 1TB serviced with 250 JVMs of 4GB each? Of course you could, but because of your organization's legacy rules from the old 32-bit JVM days, you continue to spin JVMs that are less than 1GB. More realistically, though, the notion that larger JVMs may have larger GC pauses is ill conceived. That belief is not entirely true, but not completely false either. Yes, larger pauses will occur, but with recent advancements in 64-bit JVM and concurrent mark-sweep (CMS) GC, the days of larger and less-pause-sensitive JVMs have arrived. Not only has GC gotten a lot better, but also the underlying server hardware has gotten better to support 4GB heap spaces. In fact, 4GB is a unique and magical number because JVMs these days automatically treat the 4GB heap space as 32-bit address space within a 64-bit JVM to save on memory usage. This is possible because a 32-bit address range is within 4GB. In fact, the JVM using the `-XX:+UseCompressedOops` option can be applied to Java heap spaces of up to 32GB.

After reading this book, you will understand that alternatives and workloads suited for larger JVMs exist. Clearly, I am not advocating larger JVMs for everyone, but a 4GB JVM is really not that big anymore. Keep in mind that I *am* advocating a more reasonable number of JVMs, even if that means increasing the heap size. And remember, if you have a vendor that says that you will incur a performance cost when moving from 32-bit to 64-bit JVMs, this is not entirely true. Our compression optimization experiences have mostly disproved the notion that migrating from a 32-bit to a 64-bit JVM causes performance degradation. Consider, for instance, servicing 1TB in 250 JVMs versus 1000 JVMs. Ask the vendor how much you save by not having to run 750 JVMs (because you went from a 1GB JVM heap to a 4GB JVM heap). The cost savings would include some of the 750 GC cycles that you no longer use in addition to the underlying CPU cores that you free up. You also save because you do not have to pay for additional licenses.

The latter chapters in this book delve into various sizes of JVMs and when you would use one versus the other.

## Why Should I Bother with Virtualization? What Are Some Key Benefits?

Perhaps 5 years ago, we still had some customers asking the “why virtualize” question. In recent years, though, the benefits of virtualization have become widely understood as virtualization has pretty much become the standard. This standard is based on VMware virtualization technology, mostly because of its robustness and its fifth-generation maturity.

Virtualization offers the following key benefits:

- **Mature, proven, and comprehensive platform:** VMware vSphere (<http://www.vmware.com/products/datacenter-virtualization/vsphere/overview.html>) is fifth-generation virtualization (many years ahead of any alternative). It delivers higher reliability, more advanced capabilities, and greater performance than competing solutions.
- **High application availability:** High-availability infrastructure remains complex and expensive. But VMware integrates robust availability and fault tolerance right into the platform to protect virtualized applications. Should a node or server ever fail, all the VMs are automatically restarted on another machine.
- **Wizard-based guides for ease of installation:** VMware’s wizard-based guides take the complexity out of setup and configuration. You can be up and running in one-third the deployment time of other solutions.
- **Simple, streamlined management:** VMware lets you administer both your virtual and physical environments from a “single pane of glass” console right on your web browser. Time-saving features such as auto-deploy, dynamic patching, and live VM migration reduce routine tasks from hours to minutes. Management becomes much faster and easier, boosting productivity without adding to your headcount.
- **Higher reliability and performance:** Our platform blends CPU and memory innovations with a compact, purpose-built hypervisor that eliminates the frequent patching, maintenance, and I/O bottlenecks of other platforms. The net result is best-in-class reliability and consistently higher performance (for heavy workloads, two-to-one and three-to-one performance advantages over our nearest competitors).
- **Superior security:** VMware’s hypervisor is much thinner than any rival, consuming just 144MB compared with others’ 3 to 10GB disk profile. Our small hypervisor footprint presents a tiny, well-guarded attack surface to external threats, for airtight security and much lower intrusion risk.
- **Greater savings:** VMware trumps other virtualization solutions by providing 50% to 70% higher VM density per host—elevating per-server utilization rates from 15% to as high as 80%. You can run many more applications on much less hardware than with other platforms, for significantly greater savings in capital and operating costs.

- **Affordability:** VMware is highest in capabilities, but not cost. Starting at \$165 per server, the small business packages consolidate more of your applications on fewer servers, with greater performance—delivering the industry’s lowest total cost of ownership (TCO).

To quickly determine and compare the cost of deploying VMware virtualization in your environment, use the VMware cost-per-application calculator at <http://www.vmware.com/go/costperappcalc>.

Fundamentally, because Java is independent of the operating system, it is the perfect candidate for virtualization because it does not have any hardware dependencies. Java also benefits from the many virtualization features such as high availability (HA) and VMotion (the ability to move VMs from one vSphere host to another without downtime). This type of agility that virtualization adds to a Java platform is fairly critical for Java platforms in general, but more specifically for large-scale Java. In large-scale Java, we often find thousands of JVMs that require constant administration and management (for instance, starting them, stopping them, and upgrading them without downtime). This type of administration activity cannot be feasibly accommodated at such a large scale without virtualization agility, such as VMotion and HA.

Enterprise Java application requirements for dynamic scalability, rapid provisioning, and HA represent a growing concern for development and operations groups today. Achieving these requirements with platforms that are completely based on conventional hardware is complex and expensive. Virtualization is a breakthrough technology that alleviates the pressures that common enterprise Java application requirements may impose on an organization. Features such as horizontal scalability, vertical scalability, rapid provisioning, enhanced HA, and business continuance are some of the key attributes available with the VMware vSphere suite. Chapter 1 examines three categories of large-scale Java platforms so that you can further appreciate the complexities of such systems.

Now that you understand that large-scale Java platforms require agility features of virtualization, and that Java operating system independence makes it a prime candidate for workload virtualization, let’s take a closer look at this Java independence from the operating system and the hypervisor.

### **Should I Virtualize Java Platforms?**

For those who don’t have time to read this entire section, I can simply answer, “Yes, you should virtualize.” After all, Java is independent of the underlying hypervisor, such as VMware’s bare-metal hypervisor, and the operating system. For those who want to delve a little more into what this means, though, read on.

The main design tenets of Java are based on a cross-platform language that is operating system independent (as long as there is an operating-system-supported underlying runtime). We know this runtime as the JVM, which has become a permanent fixture of many enterprise application platforms. You could write a Java application and run it on various JVMs on different operating systems (without having to recompile). Of course, many VMware customers have one vendor-targeted JVM in production and so would not have to worry about moving a Java application from one JVM implementation to another. If they chose to do so, however, they could easily do it, primarily because of Java's cross-platform and operating system independence facilitated by a JVM.

So, you can reasonably conclude that the Java applications do not really care which JVM is being targeted for them to run on and are independent of the specific JVM implementation and operating system.

Of course, you might ask, "What about all the different internal behaviors of one JVM versus another?" At the end of the day, they all adhere to the JVM spec, and although some JVM options (-XX, flags, and so on) have different names, they more or less behave in a similar manner. The differences are not in the language, but in the way the Java process can be optimized with various JVM options passed at the Java command line.

Now fast forward to the infrastructure side of things. VMware ESXi is a bare-metal hypervisor that makes it possible to run multiple operating systems on a particular piece of hardware. Infrastructure administrators no longer have to worry about installing one kind of operating system for one piece of hardware versus another. VMware makes the operating system run independently of the underlying hardware (bare metal) and creates a degree of independence between the operating system and the bare metal/hardware.

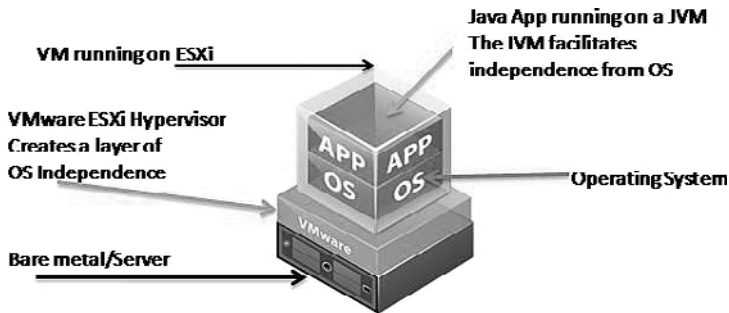
Although the answer to whether Java is both OS and hypervisor independent is clearly yes, it is due to two degrees of independence. The first degree is that of Java's main tenet of cross-platform and OS independence, and the second degree is that of VMware ESXi hypervisor making the operating system independent of the hardware that it runs on. In fact, when a Java application runs on an operating system that is in an ESXi-based VM, ESXi has no notion of whether it is a Java workload running on the operating system, making the ESX hypervisor completely independent of the workload running on it. A further testament to this is that, because of this independence, no operating system changes are needed when you deploy a Java application on a VM.

Conversely, the JVM doesn't really know that it is running on a VM sitting on an ESXi hypervisor, and to the JVM, the VM appears like any other server with compute resources (CPUs, RAM, and so on) presented to it.

As long as the JVM you are using is supported on the operating system on which your applications are running, there is no need for additional concern about or dependency on support from the downstream VM and ESXi layers.



Figure I-1 illustrates all the layers discussed in this section.



**Figure I-1** Enterprise Java Application Running on a VM Virtualized by VMware ESXi

## Who Should Read This Book?

This book is targeted at IT professionals who are in search of implementation guidelines for running enterprise Java applications on VMware vSphere in production and in QA/test environments.

The first three chapters are beneficial to CIOs, VPs, directors, and enterprise architects looking for key high-level business propositions for virtualizing enterprise Java applications. The remaining chapters are for developers and administrators looking for implementation details.

## How to Use This Book

This book consists of seven chapters, an appendix, and a glossary:

- **Chapter 1, “Introduction to Large-Scale Java Platforms”:** This chapter introduces various types of large-scale Java platforms and highlights the unique performance enhancements they require based on their scale.
- **Chapter 2, “Modern Scalable Data Platforms”:** This chapter details how modern data platforms are structured.
- **Chapter 3, “Tuning Large-Scale Java Platforms”:** This chapter highlights key considerations and provides guidelines to IT architects who are in the process of sizing their enterprise Java applications to run on VMware vSphere. This chapter explains how to obtain the best sizing configuration for your Java applications running on VMware vSphere. You are guided through the process of performance benchmarking on an application and given pointers on what to measure, what is available to be tuned, and how to best determine the optimal size for your Java application.

- **Chapter 4, “Designing and Sizing Large-Scale Java Platforms”:** This chapter walks the reader through various approaches in sizing modern virtual Java platforms. It takes the reader through actual methodology of vertical and horizontal scalability as it applies to large-scale Java platforms, while also showing actual sizing examples that can be leveraged on production systems.
- **Chapter 5, “Performance Studies”:** This chapter summarizes some of the key highlights from published performance papers.
- **Chapter 6, “Best Practices”:** This chapter provides information about best practices for deploying large-scale Java applications on VMware, including key best-practice considerations for architecture, performance, designing and sizing, and high availability. This information is intended to help IT architects successfully deploy and run Java environments on VMware vSphere.
- **Chapter 7, “Monitoring and Troubleshooting Primer”:** This chapter summarizes what to do when you hit a bottleneck or a performance issue while virtualizing Java. It provides a helpful summary for your use out in the field.
- **Appendix, “FAQs”:** This appendix is a collection of many questions from VMware customers that the author has encountered over the years. It is always helpful to quickly ramp up on any technology by reading FAQs.
- **Glossary**

*This page intentionally left blank*

## About the Author

**Emad Benjamin** has been in the IT industry for the past 20 years. He graduated with a Bachelor of Electrical Engineering degree from the University of Wollongong. Early in his career, he was a C++ software engineer. Then, in 1997, he took on his first major project using Java and has focused on Java ever since. For the past 8 years, his main concentration has been Java on VMware vSphere. Emad is a featured speaker at VMworld, SpringOne, UberConf, NFJS, and various other Java user groups around the world. Currently, Emad is a principal architect in the Global Center of Excellence focused on VMware virtualization, providing training and evangelism to various corners of the world.

## About the Technical Reviewer

**Michael Webster** is a VMware Certified Design Expert (VCDX-066) on vSphere 4 and 5, vExpert 2012–2013, and the owner of IT Solutions 2000 Ltd., which delivers project management, ITIL-based VMware operational readiness, and technical architecture consulting services to enterprise and service provider clients around the world. He has been using VMware products since 1998 and has been designing and deploying VMware solutions since 2002. He specializes in the design and implementation of virtualization solutions for Unix to Linux migrations, business-critical applications, disaster avoidance, mergers and acquisitions, and public and private cloud. Michael has been in the IT industry since 1995 and consulting since 2001. As of February 2012, IT Solutions 2000 Ltd. was granted the VMware Virtualizing Business Critical Applications (VBCA) Competency, making it one of the first companies in the world to achieve this accreditation. IT Solutions 2000 Ltd. is one of very few companies worldwide accredited to deliver projects for all the business-critical applications covered by the VBCA program: SAP, Oracle, MS SQL Server, MS Sharepoint, and MS Exchange. Michael is regularly called on to consult and speak on all aspects of virtualizing business-critical applications at events and for organizations all across the globe. Longwhiteclouds.com was recently voted one of the top 25 virtualization blogs in the world as listed on vSphere-Land.com.

# Acknowledgments

I first want to thank my wife, Christine, and our boys, Anthony and Adrian, for their understanding about not spending enough time with them while I was writing this book. Christine, you have been my pillar of strength, always understanding and accommodating.

I also want to thank my parents for sacrificing so much of their life to help me pursue my education and career, and to thank my brothers and sisters for their encouragement. I also want to thank Christine's family for their love and support.

I want to also thank my dear friend, His Grace Mar Awa Royel, Bishop of the Assyrian Church of the East, for his blessings.

I want to extend special thanks to Matt Stepanski, VP of GTS and Steve Beck, Sr. Director of GCOE, for their continuous support and encouragement with the publication of this book.

I want to extend sincere gratitude for Michael Webster's efforts in thoroughly reviewing the book. I greatly appreciate his enthusiasm and ability to promptly review the book, pointing out some key changes.

I would also like to thank my colleagues at VMware who helped make this book a reality: Lyndon Adams, Mark Achtemichuk, John Arrasjid, Scott Bajtos, Stephen Beck, Channing Benson, Jeff Buell, Dino Ciciarelli, Blake Connell, Ben Corrie, Melissa Cotton, Bhavesh Davda, Scott Deeg, Carl Eschenbach, Duncan Epping, Jonathan Fullam, Alex Fontana, Filip Hanik, Bob Goldsand, Jason Karnes, Jeremy Kuhnash, Ross Knippel, Gideon Low, Catherine Johnson, Mark Johnson, Kannan Mani, Sudhir Menon, Justin Murray, Vas Mitra, Avinash Nayak, Mahesh Rajani, Jags Ramnarayan, Raj Ramanujam, Harold Rosenberg, Dan Smoot, Randy Snyder, Lise Storc, Matt Stepanski, Mike Stolz, Guillermo Tantachuco, Don Sullivan, Abdul Wajid, Sumedh Wale, Yvonne Wassenaar, Michael Webster, Mark Wencek, James Williams, and Matthew Wood.

*This page intentionally left blank*

## **We Want to Hear from You!**

As the reader of this book, *you* are our most important critic and commentator. We value your opinion and want to know what we're doing right, what we could do better, what areas you'd like to see us publish in, and any other words of wisdom you're willing to pass our way.

We welcome your comments. You can email or write us directly to let us know what you did or didn't like about this book—as well as what we can do to make our books better.

*Please note that we cannot help you with technical problems related to the topic of this book.*

When you write, please be sure to include this book's title and author as well as your name, email address, and phone number. We will carefully review your comments and share them with the author and editors who worked on the book.

Email: [VMwarePress@vmware.com](mailto:VMwarePress@vmware.com)

Mail: VMware Press  
ATTN: Reader Feedback  
800 East 96th Street  
Indianapolis, IN 46240 USA

## **Reader Services**

Visit our website at [www.informit.com/title/9780133491203](http://www.informit.com/title/9780133491203) and register this book for convenient access to any updates, downloads, or errata that might be available for this book.



*This page intentionally left blank*

## Introduction to Large-Scale Java Platforms

This chapter defines three categories of large-scale Java platforms:

- **Category 1:** Large number of Java Virtual Machines (JVMs) (100s–1000s of JVMs)
- **Category 2:** Smaller number of JVMs with large heap sizes
- **Category 3:** A combination of category 1 consuming data from category 2

In addition, the chapter discusses various trends and outlines technical considerations to help you understand the range of technical issues associated with designing large-scale Java platforms.

### Large-Scale Java Platform Categories

Based on field interactions with customers, large-scale Java platforms typically fall into three main categories, as follows:

- **Category 1:** This category is distinguished by its large number of Java Virtual Machines (JVMs). In this category, hundreds to thousands of JVMs are deployed on the Java platform, and these are typically JVMs that function within a system that might be servicing millions of users. I have seen some customers with as many as 15,000 JVMs. Whenever you are dealing with thousands of JVM instances, you must consider the manageability cost and whether opportunities exist to consolidate the JVM instances.

- **Category 2:** This category is distinguished by a smaller number of JVMs (usually 1 to 20) but with large heap size (8GB to 256GB or higher). These JVMs usually have in-memory databases deployed on them. In this category, garbage collection (GC) tuning becomes critical, as discussed in later chapters.
- **Category 3:** The third category is a combination of the first two categories, where perhaps thousands of JVMs run enterprise applications that are consuming data from category 2 types of large JVMs in the back end.

With regard to virtualizing and tuning large-scale Java platforms, four key requirement trends hold true across these three categories:

- Compute-resource consolidation
- JVM consolidation
- Elasticity and flexibility
- Performance

Let's look at each one of these trends in more detail.

## Large-Scale Java Platform Trends and Requirements

Compute resource consolidation, JVM instance consolidation, elasticity and flexibility, and performance are some of the major trends that exist within large-scale Java platform migration projects. The following subsections examine each of these in more detail.

### Compute-Resource Consolidation

Many VMware customers find that their middleware deployments have proliferated and are becoming an administrative challenge with increasing costs. Customers, therefore, are looking to virtualization as a way of reducing the number of server instances. At the same time, customers are taking the consolidation opportunity to rationalize the number of middleware components needed to service a particular load. Middleware components most commonly run within a JVM with an observed scale of hundreds to thousands of JVM instances and provide many opportunities for JVM instance consolidation. Hence, middleware virtualization provides an opportunity to consolidate twice—once to consolidate server instances and then to consolidate JVM instances. This trend is widespread; after all, every IT shop on the planet is considering the cost savings of consolidation.

One customer in the hospitality sector went through the process of consolidating their server footprint and at the same time consolidated many smaller JVMs with a heap smaller

than 1GB. They consolidated many of these smaller 1GB JVMs into two categories: those that were 4GB and others that were 6GB. They performed the consolidation in such a manner that the net total amount of RAM available to the application was equal to the original amount of RAM, but with fewer JVM instances. They did all of this while improving performance and maintaining good service level agreements (SLAs). They also reduced the cost of administration considerably by reducing the number of JVM instances they had to originally manage; this refined environment helped them easily maintain SLAs.

Another customer, in the insurance industry, achieved the same result, but was also able to overcommit CPU in development and QA environments to save on third-party software license costs.

## **JVM Instance Consolidation**

Sometimes we come across customers that have a legitimate business requirement to maintain one JVM for an application and/or one JVM per a line of business. In these cases, you cannot really consolidate the JVM instances because doing so would cause intermixing of the lifecycle of one application from one line of business with another. However, although such customers do not benefit from eliminating additional JVM instances through JVM consolidation, they do benefit from more fully utilizing the available compute resources on the server hardware, resources that otherwise would have been underutilized in a nonvirtualized environment

## **Elasticity and Flexibility**

It is increasingly common to find applications with seasonal demands. For example, many of our customers run various marketing campaigns that drive seasonal traffic toward their application. With VMware, you can handle this kind of traffic burst by automatically provisioning new virtual machines (VMs) and middleware components when needed; you can then automatically tear down these VMs when the load subsides.

The ability to change updating/patching hardware without causing outage is paramount for middleware that supports the cloud era scale and uptime. VMware VMotion enables you to move VMs around without needing to stop applications or the VM. This flexibility alone makes virtualization of middleware worthwhile when managing large-scale middleware deployments. One customer in the financial space, handling millions of transactions per day, used VMotion quite often, without any downtime, to schedule their hardware upgrades; a process that otherwise would be costly to their business because of the required scheduled downtime.

## Performance

Customers often report improved middleware platform performance when virtualizing. Performance improvements are partly due to the updated hardware that customers will typically refresh during a virtualization project. Some performance improvement occurs, too, due to the robust VMware hypervisor. The VMware hypervisor has improved considerably in the past few years, and Chapter 5, “Performance Studies,” discusses a few performance studies done to showcase some of the heavy workloads that were tested in a virtualized environment.

## Large-Scale Java Platform Technical Considerations

When designing large-scale Java platforms, many technical considerations come into play. For example, a good understanding of Java garbage collection (GC) and of JVM architecture, hardware, and hypervisor architectures is essential to building good large-scale Java platforms. At a high level, GC, Non-Uniform Memory Architecture (NUMA), and theoretical versus practical memory limits are discussed. Later chapters provide a more detailed description, but it is imperative to start at a high-level understanding of the issues surrounding large-scale Java platform designs.

## Theoretical and Practical Limits of Java Platforms

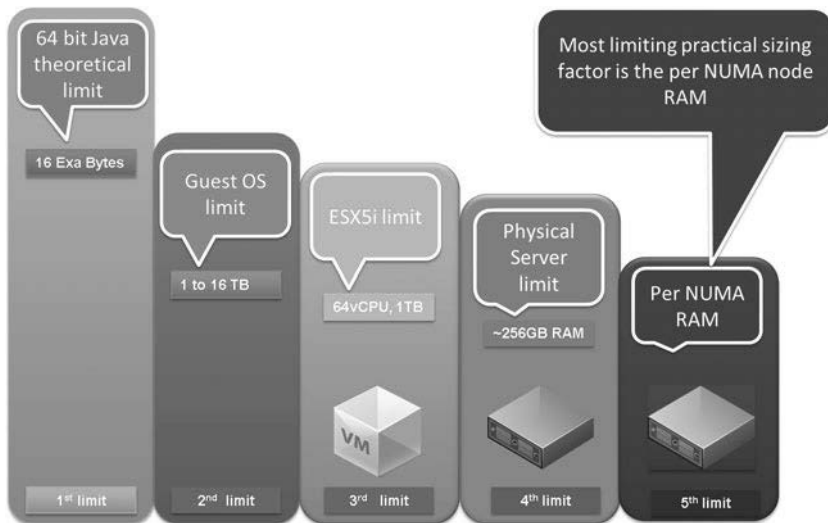
Figure 1-1 depicts the theoretical and practical sizing limits of Java workloads, critical limits to remember when sizing JVM workloads.

- It is important to highlight that the JVM theoretical limit is 16 exabytes; however, no practical system can provide this amount of memory. So, we capture this as the first theoretical limit.
- The second limit is the amount of memory a guest operating system can support; in most practical cases, this is several terabytes (TB) and depends on the operating system being used.
- The third limit is the ESXi5 1TB RAM per VM, which is ample for any workload that we have encountered with our customers.
- The fourth limit (really the first practical limit) is the amount of RAM that is cost-effective on typical ESX servers. We find that, on average, vSphere hosts have 128GB to 144GB, and at the top end 196GB to 256GB. Certainly from a feasibility standpoint, the hard limit is probably around 256GB. There are, of course, larger RAM-based vSphere hosts, such as 384GB to 1TB; however, these are probably more suited for category 2 types of in-memory database workloads and more likely

suited for traditional relational database management systems (RDBMS) that would utilize such vast compute resources. The primary reason these systems need such large vSphere hosts is because most (with some minor exceptions, such as Oracle RAC) traditional RDBMS do not scale out and mainly scale up. In the case of category 1 and category 2, a scale-out approach is available and so the potential selection of a more cost-effective vSphere host configuration is afforded. In category 1 types of Java workloads, you should consider vSphere hosts with a more reasonable RAM range of less than 128GB.

- The fifth limit is the total amount of RAM across the server and how this is divided into a number of NUMA nodes, where each processor socket will have one NUMA node worth of NUMA-local memory. The NUMA-local memory can be calculated as the total amount of RAM within the server divided by the number of processor sockets. We know that for optimal performance you should always size a VM within the NUMA node memory boundaries; no doubt, ESX has many NUMA optimizations that come into play, but it is always best to stay NUMA local.

If the ESX host, for example, has 256GB of RAM across two processor sockets (that is, it has two NUMA nodes with 128GB (256GB/2) of RAM across each NUMA node), this implies that when you are sizing a VM it should not exceed the 128GB limit for it to be NUMA local.



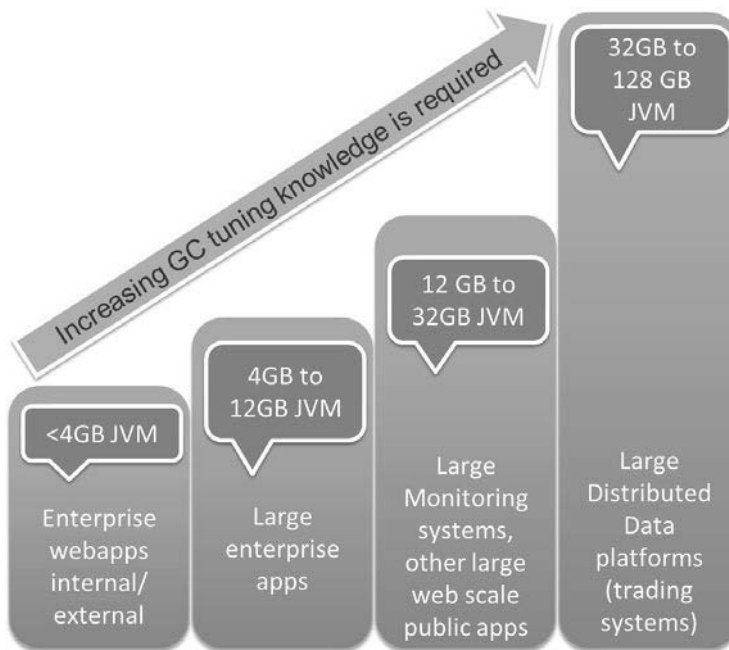
**Figure 1-1** Theoretical and Practical Limits of Java Platforms

The limits outlined in Figure 1-1 and the list will help drive your design and sizing decision as to how practical and feasible it is to size large JVMs. However, other considerations come with sizing very large JVMs, such as GC tuning complexity and knowledge needed to maintain large JVMs. In fact, most JVMs within our customer base are in the vicinity of 4GB of RAM for the typical enterprise web application, or what has been referred to in this book as category 1 workloads. However, larger JVMs exist, and we have customers that run large-scale monitoring systems and large distributed data platforms (in-memory databases) on JVMs ranging from 4GB to 128GB. This is also true for in-memory databases such as vFabric GemFire and SQLFire, where individual JVM members within a cluster can be as big as 128GB and total cluster size can be 1 to 3TB. With such large JVMs comes the need to have a better knowledge of GC tuning. At VMware, we have helped many of our customers with their GC tuning activities over the years, even though GC tuning on physical is no different from on virtual. The reason being is that we have uniquely integrated the vFabric Java and vSphere expertise into one spectrum, which has helped our customers optimally run many Java workloads on vSphere. When faced with the decision of whether to vertically scale the size of the JVM and VM, always first consider a horizontal scale-out approach; we have found that our customers get better scalability with a horizontally scaled-out platform. If horizontal scalability is not feasible, consider increasing the size of the JVM memory and hence VM memory. When opting to increase the size of the JVM by increasing the heap space/memory, the next point of consideration is GC tuning and the in-house knowledge you have to handle large JVMs.

**NOTE**

With regard to the third limit, as of this writing, ESXi 5.1 is the GA released official version; however, by the time this book is published, some of these maximum vSphere limits might change. Double-check official VMware product documentation for the latest maximums. Note, as well, that at these VM limits no cost-effective hardware would need such a large number of vCPUs; however, it is still assuring for those who might need it.

As mentioned earlier in this chapter, in the enterprise today large-scale Java platforms fall into one of three categories. Figure 1-2 shows the various workload types and relative scale. A common trend is that as the size of the JVM increases so, too, does the required JVM GC tuning knowledge.



**Figure 1-2** GC Tuning Knowledge Requirements Increase with Larger JVMs

It is important to keep the following in mind (from left to right in the figure):

- JVMs with a less than 4GB heap size are the most common among workloads today. The 4GB is a special case because it has the default advantage of using 32-bit address pointers within a 64-bit JVM space (and so has a very efficient memory footprint). These require some tuning, but not a substantial amount. This workload type falls into the realm of category 1 as defined earlier in this chapter. The default GC algorithm on server class machines is adequate. The only time you need to tune these is if the response time measurements do not suffice. In such cases, you want to follow the guidance on GC tuning in Chapter 3, “Tuning Large-Scale Java Platforms,” and in Chapter 6, “Best Practices.”
- The second workload case is still within category 1, but it is probably a serious user base internal to the organization. In this workload, we typically see heavily used (1,000 to 10,000 users) enterprise Java web applications. In these types of environments, GC tuning and slightly larger than 4GB JVMs are the norm. The DevOps team almost always has decent GC tuning knowledge and has configured the JVM away from the default GC throughput collector. Here we start to see the use of the concurrent mark and sweep (CMS) GC algorithms for these types of workloads to deliver decent response times to the user base. The CMS GC algorithm is offered



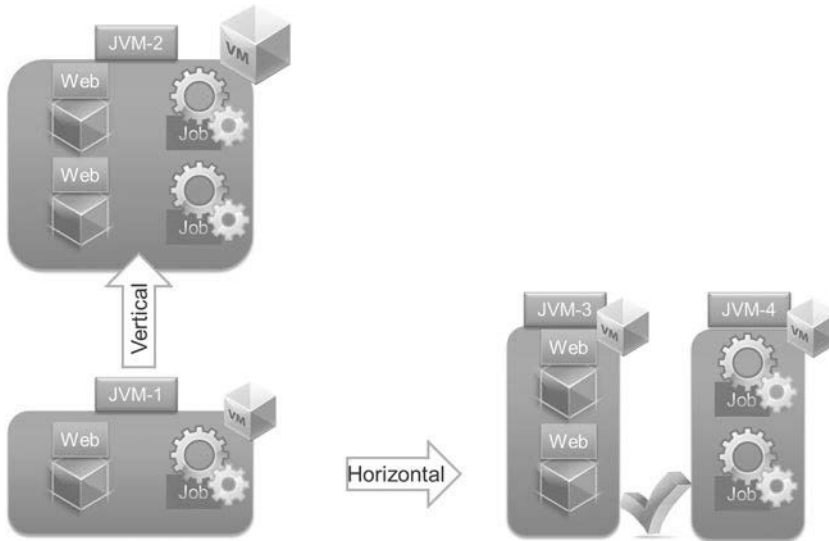
by the Oracle JVM (formerly Sun JVM). For further details and information about other GC algorithms within the Oracle JVM or IBM JVM, see Chapter 3 and Chapter 6.

- The third workload type could fall into category 2, but it is a unique case within category 2 because sometimes the larger JVMs are used because the application cannot scale out horizontally. Generic category 2 workloads are usually in-memory databases, as mentioned earlier in the chapter. In this category, a deep knowledge of JVM GC tuning is required. Your DevOps team must be able to articulate all the different GC collectors and select those most suited for improved throughput (throughput collectors) (in contrast to latency-sensitive workloads that need CMS GC to deliver better response times).
- The fourth workload type falls into both category 2 and 3. Here there could be a large distributed system, where the client enterprise Java applications are consuming data from the back-end data fabric where a handful or more of in-memory database JVM nodes are running. Tuning GC at expert level is required here.

Other than having to maintain a very large JVM, you must know the workload choices. After all, customers often scale the JVM vertically because they believe it is an easy deployment and that it is best to just leave the existing JVM process intact. Let's consider some JVM deployment and usage scenarios (perhaps something in your current environment or something you have encountered at some point):

- A customer has one JVM process deployed initially. As demand for more applications to be deployed increases, the customer does not horizontally scale out by creating a second JVM and VM. Instead, the customer takes a vertical scale-up approach. As a consequence, the existing JVM is forced to vertically scale and carry many different types of workloads with varied requirements.
- Some workloads, such as a job scheduler, require high throughput, whereas a public-facing web application requires fast response time. So, stacking these types of applications on top of each other, within one JVM, complicates the GC cycle tuning opportunity. When tuning GC for higher throughput, it is usually at the cost of decreased response time, and vice versa.
- You can achieve both higher throughput and better response time with GC tuning, but it certainly extends the GC tuning activity unnecessarily. When faced with this deployment choice, it is always best to split out the types of Java workloads into their own JVMs. One approach is to run the job scheduler type of workload in its own JVM and VM (and do the same for the web-based Java application).
- In Figure 1-3, JVM-1 is deployed on a VM that has mixed application workload types, which complicates GC tuning and scalability when attempting to scale up

this application mix in JVM-2. A better approach is to split the web application into JVM-3 and the job scheduler application into JVM-4 (that is, horizontally scaled out and with the flexibility to vertically scale if needed). If you compare the vertical scalability of JVM-3 and JVM-4 versus the vertical scalability of JVM-2 you will find JVM-3 and JVM-4 always scale better and are easier to tune.

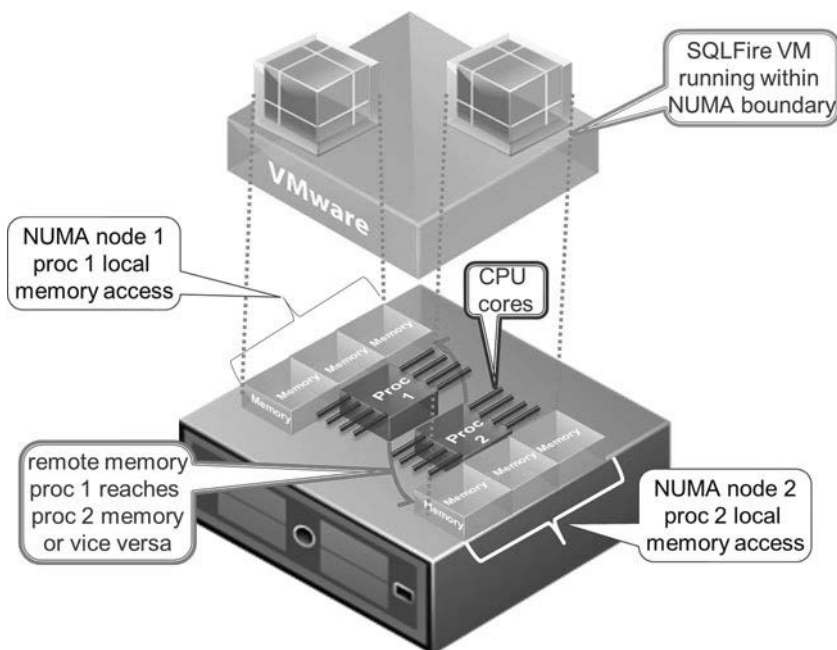


**Figure 1-3** Avoiding Mixed Workload Types in the Same JVM

## NUMA

Non-Uniform Memory Architecture (NUMA) is a computer memory design used in multiprocessors, where the memory access time depends on the memory location relative to a processor. Under NUMA, a processor can access its own local memory faster than nonlocal memory (that is, memory local to another processor or memory shared between processors).

Understanding NUMA boundaries is critical to sizing VM and JVMs. Ideally, the VM size should be confined to the NUMA boundaries. Figure 1-4 shows a vSphere host made of two sockets, and hence two NUMA nodes. The workload shown is that of two vFabric SQLFire VMs, each VM sized to fit within the NUMA node boundaries for memory and CPU. If a VM is sized to exceed the NUMA boundaries, it might possibly interleave with the other NUMA node to fulfill the request for additional memory that otherwise cannot be fulfilled by the local NUMA node. The figure depicts memory interleaving by the red arrows (dashed curved arrows show the interleaving), highlighting that this type of memory interleaving should be avoided because it may severely impact performance.



**Figure 1-4** Two-Socket Eight-Core vSphere Host with Two NUMA Nodes and One VM on Each NUMA Node

To calculate the amount of RAM available in each NUMA node, apply the equation in Formula 1-1.

$$\text{NUMA Local Memory} = \text{Total RAM on Server} / \text{Number of Sockets}$$

**Formula 1-1** Per-NUMA Node RAM Size (NUMA Local Memory)

For example, if a server has 128GB of RAM configured on it and has two sockets (as shown in Figure 1-4), this implies that the per-NUMA RAM is  $128/2$ , which equals 64GB. This is not entirely true, however, because ESX overhead needs to be accounted for. So, a more accurate approximation results from the equation shown in Formula 1-2. The formula accounts for the ESXi memory overhead (1GB as a constant, regardless of the size of the server) and a 1% VM memory overhead as 1% of the available memory. The formula is a conservative approximation, and every VM and workload will vary slightly, but the approximation should be pretty close to the worst-case scenario.

NUMA Local Memory =

$$[\text{Total RAM on Host} - \{(\text{Total RAM on Host} * \text{nVMs} * 0.01) + 1\text{GB}\}] / \text{Number of Sockets}$$

**Formula 1-2** Per-NUMA Node RAM (NUMA Local Memory) with ESXi Overhead Adjustment

The following explains the different parts of the formula:

- **NUMA Local Memory:** The local NUMA memory for best memory throughput and locality, with VM and ESXi overhead already accounted for
- **Total RAM on Host:** The amount of physical RAM configured on the physical server
- **nVMs:** The number of VMs you plan to deploy on the vSphere host
- **1GB:** The overhead needed to run ESXi
- **Number of Sockets:** The number of sockets available on the physical server, 2 socket or 4 socket

#### NOTE

Formula 1-2 assumes the most pessimistic end of the overhead range, especially as you increase the number of VMs—clearly, as you add more VMs you will have more overhead. Despite a lower number of VMs, the approximation of Formula 1-2 is pretty fair and accurate. Also, this assumes a non-overcommitted memory situation. This formula is beneficial for sizing large VMs, which is when NUMA considerations are most pertinent. When sizing large VMs, typically you are trying to maintain fewer than a handful of configured VMs, so this overhead formula accurately applies. In fact, the most optimal configuration for larger VMs that have memory-bound workloads is one VM per NUMA node. If you try to apply this formula to a deployment that has more than six VMs configured, say 10 VMs, the formula can overestimate the amount of overhead needed. More accurately, you can use the 6% rule, which maintains that regardless of the number of VMs, always assume that 6% of memory overhead is ample, whether you have 10 VMs or 20.

If you don't have time to crunch through the formula and want to quickly start configuring, assume about 6% of overhead due to memory. There are many times when not all of this is being used. For example:

**Example 1**—Using 6% approximation approach: This would imply that if you have a server which has 128GB of physical RAM (two socket hosts, eight cores on each socket) and you choose the 6% overhead approach while configuring two VMs on the host, the total NUMA local memory would be  $\Rightarrow ((128 * 0.94) - 1) / 2 \Rightarrow 59.7\text{GB}$  per VM available for memory. Because there are two VMs, the total memory offered to the two VMs is approximately  $59.7 * 2 \Rightarrow 119.32\text{GB}$ .

You also can apply the approach in Formula 1-2 as shown in Example 2 that follows:

**Example 2**—Using Formula 1-2 to calculate NUMA local available memory: Again, assuming a 128GB host with two sockets (eight cores on each socket) and two VMs to be configured on it, NUMA local memory =  $(128 - (128 * 2 * 0.01) - 1) / 2 \Rightarrow 124.44\text{GB}$ . Note that this is for two VMs. If you decide instead to configure 16 VMs of 1vCPU (1vCPU = 1 core), then the NUMA local memory per VM would be NUMA local memory =  $(128 - (128 * 16 * 0.01) - 1) / 2 \Rightarrow 53.26\text{GB}$ . This probably is overly conservative, and a more accurate representation would be around the 6% overhead calculation approach.

For best guidance, the best approximation of overhead is the 6% of total physical RAM (plus 1GB for ESXi) approach shown in Example 1.

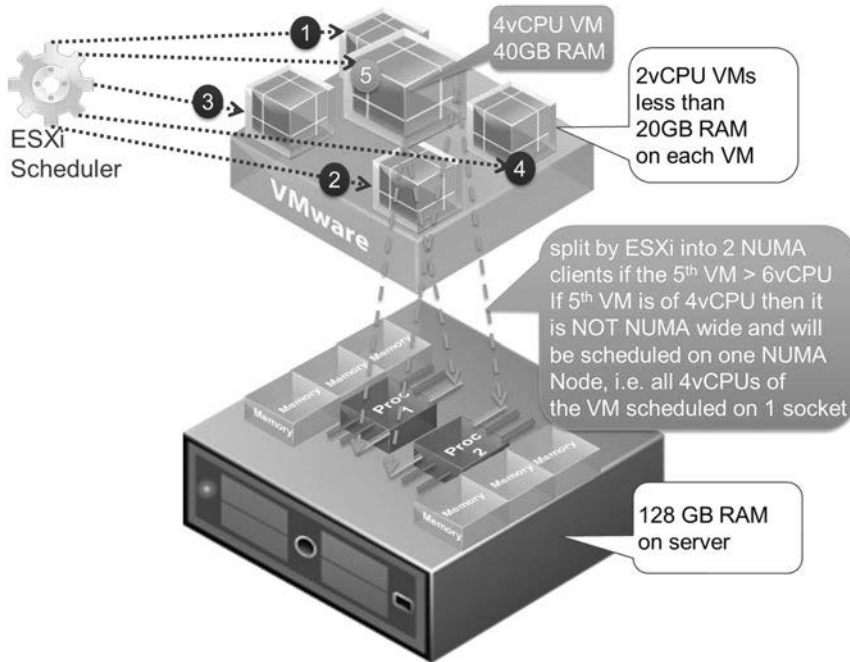
In the preceding example, showing a calculation based on a server having 128GB of RAM, the true local memory would be  $((128 * 0.99) - 1\text{GB}) / 2 \Rightarrow 62.86\text{GB}$ , which is the maximum VM size that can be configured. In this case, you can safely configure two VMs of 62.86GB of RAM and eight vCPUs each, because each of the VMs would be deployed on one NUMA node. Alternatively, you can deploy four VMs if you want to deploy smaller VMs of  $62.86\text{GB} / 2 \Rightarrow 31.43\text{GB}$  of RAM and four vCPUs each, and the NUMA scheduling algorithm would still localize the VMs to the local NUMA node.

#### NOTE

On hyperthreaded systems, VMs with a number of vCPUs greater than the number of physical cores in a NUMA node but lower than the number of logical processors (logical processors are usually shown as 2.x of physical cores, but more practically, logical processors are 1.25x of physical cores) in each physical NUMA node might benefit from using logical processors with local memory instead of full cores with remote memory. You can configure this behavior for a specific VM with the `numa.vcpu.preferHT` flag. For further details, see [http://www.vmware.com/pdf/Perf\\_Best\\_Practices\\_vSphere5.1.pdf](http://www.vmware.com/pdf/Perf_Best_Practices_vSphere5.1.pdf) and the KB article [kb.vmware.com/kb/2003582](http://kb.vmware.com/kb/2003582).

It is always advisable to start with vCPUs equal to the number of physical cores and then adjust vCPUs upward when needed, but less than approximately 1.25x of available physical cores.

To further elaborate on the ESXi NUMA scheduling algorithm, Figure 1-5 shows an example of two sockets and six cores on each socket of the server.



**Figure 1-5** ESXi NUMA Scheduling on a Two-Socket Six-Core Server

In this figure, there are initially four VMs of two vCPUs and approximately 20GB RAM on each. The initial ESXi scheduling algorithm will follow a round-robin fashion. First, step 1 occurs (as shown by the black circle with the number 1), and then the next two vCPU VMs are scheduled on the next available empty NUMA node, and then so on (steps 3 and 4) for scheduling the third and fourth VMs. At the point where all four of the 2vCPU 20GB VMs have been scheduled, and as a result of this scheduling, the four VMs will occupy the four cores on each of the sockets, as shown by red pins in this figure (red pins are the pins the four 2 vCPU VMs were initially scheduled ESXi). Moments later, a fifth VM made of four vCPUs and 40GB RAM is deployed, and now ESXi attempts to schedule this VM across one NUMA node. This is because the VM is 4vCPU and is not considered a NUMA-wide VM, so all four of its vCPUs will be scheduled on one NUMA node, even though only two vCPUs are available. What will likely happen in terms of the NUMA balancing awareness algorithm is that the ESXi scheduler will eventually force one of the two vCPU VMs to migrate to the other NUMA node in favor of trying to fit the fifth 4vCPU VM into one NUMA node. The ESXi scheduler behaves like this because it uses a concept of NUMA client and schedules VMs per NUMA client, where the default size of the NUMA client is the size of the physical NUMA node. In this case, the default is 6, so any VM that is 6vCPU and less will be scheduled on one NUMA node because it fits into one NUMA client. If you want to change this behavior, you would have to force

the NUMA client calculation to something more granular. The NUMA client calculation is controlled by *numa.vcpu.maxPerClient*, which can be set as Advanced Host Attributes -> Advanced Virtual NUMA Attributes, and if you were to change this to 2, then effectively every socket in our example will have three NUMA clients, so each 2vCPU VM will be scheduled into one NUMA client, and the fifth 4vCPU VM will be scheduled across two NUMA clients, and potentially across two sockets if need be. You seldom need get to this level of tuning, but this example illustrates the power of the NUMA algorithm within vSphere, which far exceeds any nonvirtualized Java platforms.

In general, when a virtual machine is powered on, ESXi assigns it a home node as part of its initial placement algorithm. A virtual machine runs only on processors within its home node, and its newly allocated memory comes from the home node as well. Unless a virtual machine's home node changes, it uses only local memory, avoiding the performance penalties associated with remote memory accesses to other NUMA nodes. When a virtual machine is powered on, it is assigned an initial home node so that the overall CPU and memory load among NUMA nodes remains balanced. Because internode latencies in a large NUMA system can vary greatly, ESXi determines these internode latencies at boot time and uses the information when initially placing virtual machines that are wider than a single NUMA node. These wide virtual machines are placed on NUMA nodes that are close to each other for lowest memory access latencies. Initial placement-only approaches are usually sufficient for systems that run only a single workload, such as a benchmarking configuration that remains unchanged as long as the system is running. However, this approach cannot guarantee good performance and fairness for a datacenter-class system that supports changing workloads. Therefore, in addition to initial placement, ESXi 5.0 does dynamic migration of virtual CPUs and memory between NUMA nodes for improving CPU balance and increasing memory locality. ESXi combines the traditional initial placement approach with a dynamic rebalancing algorithm. Periodically (every two seconds by default), the system examines the loads of the various nodes and determines whether it should rebalance the load by moving a virtual machine from one node to another.

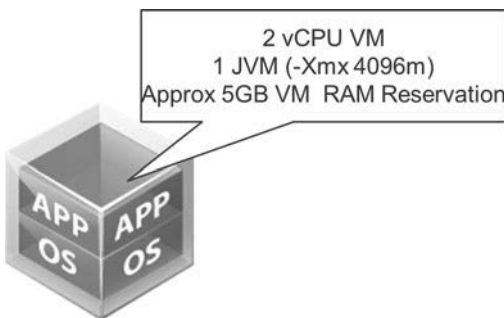
This calculation takes into account the resource settings for virtual machines and resource pools to improve performance without violating fairness or resource entitlements. The rebalancer selects an appropriate virtual machine and changes its home node to the least loaded node. When it can, the rebalancer moves a virtual machine that already has some memory on the destination node. From that point on, the virtual machine allocates memory on its new home node and runs only on processors in the new home node. Rebalancing is an effective solution to maintain fairness and ensure that all nodes are fully used. The rebalancer might need to move a virtual machine to a node on which it has allocated little or no memory. In this case, the virtual machine incurs a performance penalty associated with a large number of remote memory accesses. ESXi can eliminate this penalty by transparently migrating memory from the virtual machine's original node to its new home node.

**NOTE**

In vSphere 4.1/ESXi 4.1, the underlying physical NUMA architecture is not exposed by the hypervisor to the operating system, and therefore application workloads running on such VMs cannot take specific advantage of additional NUMA hooks that they may provide. However, in vSphere5, the concept of vNUMA was introduced, where through configuration you can expose the underlying NUMA architecture to the operating system, and so NUMA-aware applications can take advantage of it. In Java, the `-XX:+UseNUMA` JVM option is available; however, it is compatible only with the throughput GC and not the CMS GC. Paradoxically, in most memory-intensive cases where NUMA is a huge factor, latency sensitivity is a big consideration, and therefore the CMS collector is more suitable. This implies that you cannot use CMS and the `-XX:+UseNUMA` option together. The good news is that vSphere NUMA algorithms are usually good enough to provide locality, especially if you have followed good NUMA sizing best practices—such as sizing VMs to fit within NUMA boundaries for memory and vCPU perspective.

**Most Common JVM Size Found in Production Environments**

Having discussed thus far the various JVM sizes that you can deploy (in some cases, very large JVMs), it is important to keep in mind that the most common JVMs found in data centers are of 4GB heap size. This may be a fairly busy JVM with 100 to 250 concurrent threads (actual thread count will vary because it depends on the nature of the workload), 4GB of heap, approximately 4.5GB for the JVM process, 0.5GB for the guest operating system, and so a total recommended memory reservation for the VM of 5GB with two vCPUs and one JVM process, as shown in Figure 1-6.

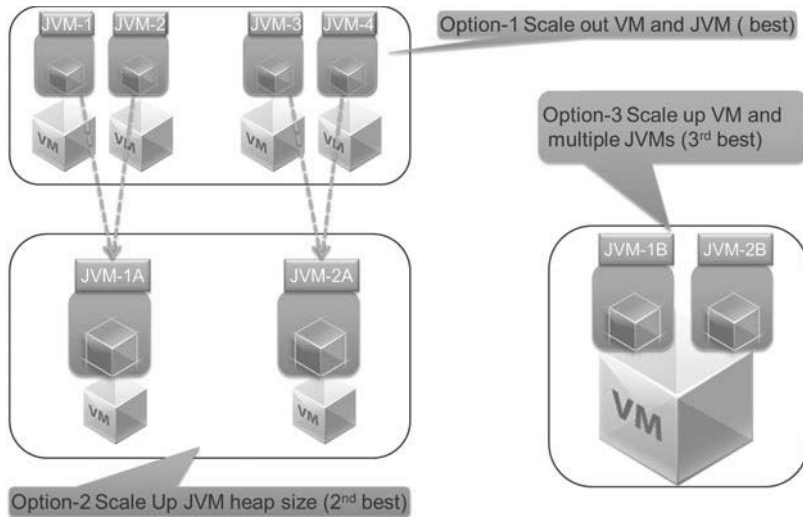


**Figure 1-6** Most Common JVM Size Found in Production Environments



## Horizontal Scaling Versus Vertical Scaling of JVMs and VMs

When considering horizontal scaling versus vertical scaling, you have three options, as Figure 1-7 shows.



**Figure 1-7** Horizontal Versus Vertical JVM Scalability Choices

The sections that follow detail the pros and cons of these three options.

### Option 1

With Option 1, JVMs are introduced into the Java platform by creating a new VM and deploying a new JVM on it (hence, a scale-out VM and JVM model).

#### Option 1 Pros

This option provides the best scalability because the VM and the JVM are scheduled out as one unit by the ESXi scheduler. It is really the VM that is scheduled by the ESXi, but because there is only one JVM on this VM, the net effect is that the VM and the JVM are scheduled as one unit.

This option also offers the best flexibility to shut down any VM and JVM in isolation without impacting the rest of the Java platform. This is no doubt in relative terms, however, because most Java platforms are horizontally scalable, and in most cases there are enough instances to service traffic, even though JVM instances are being shut down. The relative comparison in terms of more instances having better scalability is based on having 100 JVMs and VMs versus having 150 JVMs and VMs for the exact same system, if for a

specific instance you where comparing and contrasting platform design and were trying to choose between 100 JVM systems versus 150 JVMs, with both cases of 100 and 150 JVMs having the same net RAM. Clearly, the system with 150 JVMs will have the better flexibility and scalability. In the 150 JVM scenario, because you have more JVMs, it is likely that the size of the JVM is smaller compared to a system that has 100 JVMs. In this case, if a JVM from the 150 JVM platform encounters a problem, likely the impact is smaller because the JVM holds less data than in the 100 JVM scenario. So, the scale-out robustness of the 150 JVMs will prove to be more prudent.

If the system has been refined, the horizontal scalability advantages assumed previously apply. *Refined* here means that VM and JVM best practices have been applied based on a 64-bit JVM architecture having a reasonable-size JVM with an approximate minimum of 4GB heap space, and not fragmented around a legacy 32-bit JVM limit of 1GB heap space. (Some legacy 32-bit JVMs could withstand greater than 1GB, but for practical use, 32-bit JVMs have a legacy 1GB limit.)

### Option 1 Cons

This option is expensive because it leads to having more operating system copies, and licensing becomes expensive quite quickly. Administering such a system is more expensive because there are more VMs and JVMs to keep track of.

No technical reason requires you to place one JVM on one VM. The only exception is in the case of systems that are in-memory databases (like category 2) that require high throughput memory from the local NUMA node. In those cases, the VMs are sized to fit within the NUMA node and will have only one JVM on them. Also note that the JVMs in in-memory databases tend to be quite large, sometimes as big as 128GB, as opposed to category 1 JVM sizes (typically 1 to 4GB heap size). In cases such as option 1, however, which are essentially in category 1 (as defined earlier in this chapter), you have many opportunities to consolidate the JVMs and eliminate wasteful JVMs and VM instances.

This is a common pattern among legacy 32-bit JVMs, where the 1GB limit of the 32-bit JVM would have forced Java platform engineers to install more JVM instances to deal with increases in traffic. The downside here is that you are paying for additional CPU/licenses. If you consolidate JVMs by migrating to 64-bit JVMs and increasing the heap at the same time, you will save by having fewer JVMs servicing the same amount of traffic. Of course, the JVM size will likely increase from, for example, 1GB to 4GB.

### Option 2

Option 2 involves scaling up the JVM heap size by consolidating fragmented smaller JVMs and also as a result consolidate VMs.

## Option 2 Pros

The pros for using Option 2 are as follows:

- Reduced administration cost due to the lower number of JVMs and VMs
- Reduced licensing cost due to fewer operating system copies
- Improved response times as more transactions are (most likely) now executed within the same heap spaces, as opposed to requiring marshaling across the network to other JVMs
- Reduced hardware cost

### NOTE

If you look at option 2 in Figure 1-7, it shows that two JVMs (JVM-1A and JVM-2A) were consolidated from the four JVMs (JVM-1, 2, 3, and 4) of option 1. In this process, the four VMs were also consolidated into two VMs, as shown in the figure. For example, if JVM-1, 2, 3, and 4 were all of 2GB heap size, each running on VMs of 2vCPU, this implies the total RAM serviced to the heap, and in turn to the application, is 8GB across all the JVMs. The total vCPUs across all the VMs is eight vCPUs. Now when consolidating down to two VMs and two JVMs, the JVMs in option 2 (JVM-1A and JVM-2A) are each of 4GB heap, for a total of 8GB, and the VMs are two vCPUs each. This implies a total of four vCPUs across both VMs, a savings of four vCPUs, because originally in option 1 there were four VMs of two vCPUs each.

It is possible to scale down vCPUs while still maintaining an equal amount of RAM (Java heap space) because with larger JVM heap spaces, GC can scale vertically fairly well without having to excessively consume CPU. This is largely workload behavior dependent, and some workloads may indeed exhibit increased CPU usage when JVMs are scaled up. However, most category 1 workloads have exhibited a behavior of releasing the unneeded vCPU when consolidated into a larger JVM heap. 64-bit JVMs are highly capable runtime containers, and although there is an initial cost of launching one, they do enable you to crunch through a massive number of transactions that are within much larger heap spaces. When you are thinking about creating a new JVM, you want to ask the same questions as if you were about to create a new VM. If someone needs a new VM, a vSphere administrator always asks why it is needed. Because the JVM is a highly capable machine (just as the VM is a highly capable compute resource), vSphere administrators and DevOps engineers should always scrutinize whether the creation of a new JVM is necessary (as opposed to leveraging existing JVM instances and perhaps increasing the heap space, within reason, to facilitate more traffic).

## Option 2 Cons

The cons for using Option 2 are as follows:

- Because of the larger-size JVMs, you risk losing more data (when compared with the case of smaller JVMs in option 1) if a JVM crashes without proper redundancy or persistence of transactions in place.
- Due to consolidation, you might have fewer high-availability (HA) JVM instances.
- Consolidation is limited to line of business. You do not want to mix applications from different lines of business into the same JVM; a crash of the JVM would impact both lines of business if you were to mix two into one JVM.
- Larger JVMs may require some more GC tuning.

## Option 3

If option 1 and option 2 are not possible, consider option 3. In this case, you are placing multiple JVMs on a larger VM. Now JVM-1B and JVM-2B could be JVMs that are consolidated copies, like the ones in option 2, or nonconsolidated copies like in option 1. In either case, you can stack these JVMs on a larger VM, or multiple large VMs for that matter.

## Option 3 Pros

The pros for using Option 3 are as follows:

- If the current platform is similar to that in option 1, it might be an advantage, because of logistical reasons, to keep the current number of JVMs intact in the deployment, but then consider building larger VMs with multiple JVMs stacked on them.
- Reduced number of operating system licenses.
- Reduced number of VM instances.
- Reduced administration cost due to having fewer VMs.
- You can have dedicated JVMs for each line of business but can also deploy JVMs from multiple lines of business on the same VM. You should do this only if the cost of VM consolidation outweighs the danger of having multiple lines of business impacted during a VM crash.
- Large VMs makes it possible to have more vCPUs for JVMs. If a VM has two large JVMs on it from different lines of business, for example, and they peak at different times, it is likely that all the vCPUs are available to the busy JVM, and then similarly for the next JVM when its peak arrives.

### Option 3 Cons

Larger VMs will most likely be required. Scheduling larger VMs may require more tuning than smaller VMs.

#### NOTE

Various performance studies have shown that the sweet-spot VM size is two vCPUs to four vCPUs for category 1 workloads. Category 2 workloads require more than four vCPUs; at a minimum, four vCPUs may be needed. Remember, though, that scheduling opportunity from an HA perspective may be diminished. However, category 2 workloads, as in-memory databases, are mostly fault tolerant, redundant, and disk persistent, and therefore might not rely as much on VMware HA or automatic Distributed Resource Scheduler (DRS).

Because this option is about trying to consolidate VMs, it is highly likely that JVMs from different lines of business may be deployed on the same VM. You must manage this correctly because inadvertent restart of a VM may potentially impact multiple lines of business.

You can attempt to consolidate JVMs in this case and also stack them up on the same VM; however, this forces the JVMs to be much larger to fully utilize the underlying memory. If you configure fewer larger VMs, it literally means that you have VMs with a lot more RAM from the underlying hardware, and to fully consume this you might need larger JVM heap spaces. Because of the larger-size JVMs, you risk losing more data if a JVM crashes without proper redundancy or persistence of transactions in place.

This option might require large vSphere hosts, and larger servers cost more.

## Summary

This chapter introduced the concept of large-scale Java platforms and described how they generally fall into one of three categories:

- **Category 1:** Large number of JVMs
- **Category 2:** Smaller number of JVMs with large heap sizes
- **Category 3:** A combination of category 1 and 2

The chapter also examined the various theoretical and practical limits that exists within the JVM and outlined various workload types and commonly encountered JVM sizes. The chapter also discussed the NUMA and the various pros and cons of horizontal scalability, vertical scalability, JVM consolidation, and VM consolidation.

## A

- accessing, performance charts, 195
- active-active architectures (SQLFire), 57, 60
- active memory counter (vCenter), 192
- adjusting young generation internal cycle, 72-74
- after-events, 47
- AlwaysPreTouch JVM configuration option, 69
- AMQP (Advanced Messaging Queuing Protocol), 216
- Apache Derby project, 22
- application servers
  - IBM WebLogic, virtualizing, 212
  - IBM WebSphere Application, 221
  - ratio of server threads to database, 95
  - recommendations for vSphere, 219
- asynchronous listeners (SQLFire), 32, 52-54
- available RAM per NUMA node, calculating, 10-12

## B

- benchmarks
  - environment, sizing, 91, 95
  - establishing for horizontal scalability test, 86-91
  - establishing for vertical scalability test, 82-83, 86
- best practices
  - category 1
    - high-level vSphere, 165-166
    - horizontal scalability, 158-160
    - inter-tier configuration, 160-164
    - vCPU for VM, 149-150
    - vertical scalability, 156-157
    - VM memory size, 150-155
    - VM sizing and configuration, 148
    - VM timekeeping, 156
  - category 2, 166-167
    - SQLFire, 168-172
    - SQLFire on vSphere, 173-181
  - category 3, IBM JVM and Oracle Rokit JVMs, 181-183
    - global multisite topology (SQLFire), 29-30
    - peer-to-peer topology (SQLFire), 27
    - server groups (SQLFire), 34
    - SQLFire client/server topology, 26
- building block VM, 83, 105-106, 148

## C

- calculating
  - available RAM per NUMA node, 10-12
  - number of vFabric SQLFire members, 105-106
  - VM memory requirements, 84-86
- category 1 workloads, 145
  - best practices
    - high-level vSphere, 165-166
    - horizontal scalability, 158-160
    - inter-tier, 160-164
    - vCPU for VM, 149-150
    - vertical scalability, 156-157
    - VM memory size, 150-155
    - VM sizing and, 148
    - VM timekeeping, 156
- category 2 workloads, 27, 145
  - best practices, 166-167
    - SQLFire, 168-172
    - SQLFire on vSphere, 173-181
  - sizing SQLFire Java platforms
    - entities, 97-100
    - memory, 100, 104-105
    - number, 105-106
- category 3 workloads, 145
  - best practices, IBM JVM and Oracle, 181-183
  - client/server topology (SQLFire), 25
- client/server topology (SQLFire), 24
  - best practices, 26
  - client tier, 25
  - server tier, 25
- client tier (SQLFire client/server topology), 25
- clusters, 147-148, 158-160
- CMS (concurrent mark sweep), JVM configuration
  - options, 64-66
- CMSInitiatingOccupancyFraction=75 JVM
  - configuration option, 67
- collecting metrics from vCenter, 191-193
- colocation (SQLFire), 38-39
- comparing SQLFire performance versus RDBMS, 121-123
  - CPU utilization, 125
  - response time, 124
  - scalability, 124

compute-resource consolidation, 2-3  
consolidation, 208  
consumed memory counter (vCenter), 192  
cookies, troubleshooting SYN cookies, 179  
CPU  
    HT, 83, 224  
    thread ratios, establishing, 93  
    utilization on SQLFire, comparing with RDBMS, 125  
    vCenter performance charts, 193  
    vCPU for VM best practices, 149-150

## D

daily transactional data, selecting, 97  
data fabric, identifying data size for category 2 workloads, 100, 104-105  
DB server tier, 190  
DBSynchronizer, 32  
DBSynchronizer (SQLFire), 54-57  
DDUtils, 32  
deploying new environments  
    benchmark, establishing, 82-91  
    production environment, sizing, 95  
    workload profile, establishing, 80-81  
DisableExplicitGC JVM configuration option, 69  
disk charts (vCenter), 194  
disk persistence (SQLFire), 39-41  
downloading Sprint Travel, 122  
DRS (VMware Distributed Resource Scheduler), 218

## E

EJBs (Enterprise JavaBeans), 205  
elasticity in large-scale Java platforms, 3  
entity groups, 97  
    identifying for category 2 workloads, 97-100  
establishing  
    benchmarks  
        for horizontal scalability test, 86-91  
        for vertical scalability test, 82-83, 86  
    building block VM, 105-106  
    thread ratios, 93  
    workload profile, 80-81  
ESXi 3, comparing performance with ESXi 4.1 and ESXi 5, 139-141  
ESXi 4.1, comparing performance with ESXi 3 and ESXi 5, 139-141  
ESXi 5, 6  
    comparing performance with ESXi 3 and ESXi 4, 139-141  
    NUMA scheduling algorithm, 13-14  
esxtop, troubleshooting techniques, 195, 198  
example of SQLFire sizing, 112, 115-119

## F

FAQs, 206-228  
fast data, 112  
features of SQLFire, 22  
    asynchronous listeners, 52-54  
    colocation, 38-39  
    DBSynchronizer, 54-56  
    DDLUtils, 57  
    disk persistence, 39-41  
    listeners, 47-49  
    partitioning, 31, 34-37  
    redundancy, 31-32, 37-38  
    RowLoader, 46-47  
    server groups, 30-34  
    transactions, 41-46  
    writers, 50-51  
flexibility in large-scale Java platforms, 3

## G

GC (garbage collection)  
    CMS/parallel GC configuration options, 64-66  
        AlwaysPreTouch JVM configuration option, 69  
        CMSInitiatingOccupancyFraction=75 JVM configuration option, 67  
        DisableExplicitGC JVM configuration option, 69  
        MaxTenuringThreshold=15 JVM configuration option, 68  
        OptimizeStringConcat JVM configuration option, 69  
        ParallelGCThreads JVM configuration option, 68-69  
        ScavengeBeforeFullGC JVM configuration option, 67  
        SurvivorRatio JVM configuration option, 67  
        TargetSurvivorRatio=80 JVM configuration option, 67  
        UseBiasedLocking JVM configuration option, 67  
        UseCMSInitiatingOccupancyOnly JVM configuration option, 67  
        UseCompressedOops JVM configuration option, 69  
        UseCompressedStrings JVM configuration option, 69  
        UseConcMarkSweepGC JVM configuration option, 66  
        UseNUMA JVM configuration option, 69  
        UseParNewGC JVM configuration option, 67  
        UseStringCache JVM configuration option, 69  
        -Xmn21g JVM configuration option, 66

- policy selection, 184
  - for IBM JVM, 186
  - for Oracle HotSpot JVM, 184-185
  - for Oracle jRockit, 187
- throughput GC, 63-64
- tuning
  - old generation tuning, 76
  - survivor spaces tuning, 78
  - young generation tuning, 71-74
- GemFire, 22
- global multisite topology (SQLFire), 28
  - best practices, 29-30
- granted memory counter (vCenter), 192

## H

- HA (high availability), 218
- heap, off-heap section, 84, 107
- high-level vSphere best practices, 165-166
- horizontal scalability
  - benchmark for testing, establishing, 86-91
  - best practices, 158-160
  - of JVMs and VMs, 16-20
  - vSphere features, 218
- host clusters, best practices, 158-160
- host CPU utilization, 224
- HotSpot, GC policy selection, 184-185
- HotSpot JVM, internal memory sections, 106
- HT (hyperthreading), 12, 83, 224

## I

- IBM JVM
  - best practices, 181-183
  - GC policy selection, 186
- IBM WebLogic application servers, virtualizing, 212
- IBM WebSphere
  - customer references, 222
  - licensing, 221
  - support for, 221
- identifying
  - entity groups for category 2 workloads, 97-100
  - memory size of data fabric for category 2, 100, 104-105
  - number of required SQLFire members, 105-106
- in-memory data management systems, SQLFire, 22
  - active-active architectures, 57, 60
  - Apache Derby project, 23
  - asynchronous listeners, 52-54
  - client/server topology, 24-26
  - colocation, 38-39
  - DBSynchronizer, 54-56
  - DDLUtils, 57
  - disk persistence, 39-41
  - enterprise data fabric system, 24

- features, 22, 32
- global multisite topology, 28-30
- listeners, 47-49
- low-end volume, 24
- partitioning, 31, 34-37
- peer-to-peer topology, 27
- redundancy, 31, 37-38
- redundancy zones, 28
- RowLoader, 46-47
- server groups, 30-34
- share-nothing persistence mechanism, 22
- transactions, 41-46
- writers, 50-51

- inspecting thread dumps, 204

- internal memory sections of HotSpot JVM, 106

- inter-tier configuration best practices, 160-164

## J-K

- Java application server tier, 190

- Java platforms

- coding practices on vSphere, 217
- compute-resource consolidation, 2-3
- elasticity, 3
- flexibility, 3
- JVM instance consolidation, 3
- memory, troubleshooting, 202-203
- NUMA, 9, 12
- technical considerations, 4
- theoretical and practical limits of, 4-6
- thread contentions, troubleshooting, 203-204
- troubleshooting, 198, 201
- tuning knowledge requirements, 7-9

- JConsole, troubleshooting Java, 198-201

- JVM Perm Size, 107

- JVMs (Java virtual machines)

- compute-resource consolidation, 2-3
- configuration options
  - AlwaysPreTouch, 69
  - CMSInitiatingOccupancyFraction=75, 67
  - DisableExplicitGC, 69
  - MaxTenuringThreshold=15, 68
  - OptimizeStringConcat, 69
  - ParallelGCThreads, 68-69
  - ScavengeBeforeFullGC, 67
  - SurvivorRatio, 67
  - TargetSurvivorRatio=80, 67
  - UseBiasedLocking, 67
  - UseCMSInitiatingOccupancyOnly, 67
  - UseCompressedOops, 69
  - UseCompressedStrings, 69
  - UseConcMarkSweepGC, 66
  - UseNUMA, 69
  - UseParNewGC, 67



UseStringCache, 69

-Xmn21g, 66

horizontal versus vertical scaling, 16-20

HotSpot JVMs, internal memory sections, 106  
IBM

best practices, 181-183

GC policy selection, 186

policy selection, 186

instance consolidation, 3

most common size in production environments,  
15

Oracle HotSpot, GC policy selection, 184-185

Oracle Rockit, best practices, 181-183

sizing, 108-111, 227-228

stacking, 83

## L

large JVMs, NUMA implications for sizing, 108-111  
licensing, 221

listeners (SQLFire), 47-49

load balancer tier, 190

load profile

establishing, 80-81

properties, 80-81

loaders, RowLoader, 46

logical processors on hyperthreaded systems, 12

## M

MaxTenuringThreshold=15 JVM configuration  
option, 68

memory

HotSpot JVM, internal memory sections, 106  
interleaving, 10

Java, troubleshooting Java, 202-203

memory size of data fabric, determining for  
category, 100, 104-105

NUMA, 9, 12

available RAM per node, calculating, 10-12

ESXi NUMA scheduling, 13-14

large VMs and JVMs, sizing, 108-111

vNUMA, 15, 110

off-the-heap area, 107

VM memory size best practices, 150-155

VM requirements, calculating, 84-86

memory counters (vCenter), 192

metrics

collecting from vCenter, 191-193

esxtop, 195, 198

migration

phases for tiers, 220

SMA, 205-207

monitoring system, tuning, 225-226

multipass collector, 64

## N

network performance charts (vCenter), 194

new environments, deploying

benchmark, establishing, 82-83, 86-91

production environment, sizing, 95

workload profile, establishing, 80-81

NUMA (Non-Uniform Memory Architecture), 9, 12,  
139-140

available RAM per node, calculating, 10-12

ESXi NUMA scheduling algorithm, 13-14

large VMs and JVMs, sizing, 108-111

vNUMA, 15, 110

NUMA node interleave, 105

## O

off-heap section, 84

off-the-heap area, 107

old generation tuning (GC), 76

Olio performance study, 127, 131

opening support-request tickets, 191

OptimizeStringConcat JVM configuration option,  
69

Oracle Hotspot JVM, GC policy selection, 184-185

Oracle jRockit JVMs

best practices, 181-183

GC policy selection, 187

Oracle WebLogic. *See* WebLogic

OutOfMemory error (Java), troubleshooting, 202

## P

parallel GC

AlwaysPreTouch JVM configuration option, 69

CMSInitiatingOccupancyFraction=75 JVM  
configuration option, 67

DisableExplicitGC JVM configuration option, 69

MaxTenuringThreshold=15 JVM configuration  
option, 68

OptimizeStringConcat JVM configuration  
option, 69

ParallelGCThreads JVM configuration option,  
68-69

ScavengeBeforeFullGC JVM configuration  
option, 67

SurvivorRatio JVM configuration option, 67

TargetSurvivorRatio=80 JVM configuration  
option, 67

UseBiasedLocking JVM configuration option, 67

UseCMSInitiatingOccupancyOnly JVM

configuration option, 67

UseCompressedOops JVM configuration option,  
69

UseCompressedStrings JVM configuration

option, 69

- UseConcMarkSweepGC JVM configuration option, 66
- UseNUMA JVM configuration option, 69
- UseParNewGC JVM configuration option, 67
- UseStringCache JVM configuration option, 69
- Xmn21g JVM configuration option, 66
- ParallelGCThreads JVM configuration option, 68-69
- partitioning (SQLFire), 31, 34-37, 99
- peer-to-peer topology (SQLFire), 27
- performance
  - ESXi 3, comparing with ESXi 4.1 and ESXi 5, 139-141
  - GC tuning
    - old generation tuning, 76
    - survivor spaces tuning, 78
    - young generation tuning, 71-74
  - Java applications on vSphere, 207
  - large-scale Java platforms, 4, 7-9
  - NUMA node interleave, 105
  - Olio performance study, 127, 131
  - SpringTrader performance study, 131-132
    - application and, 133-137
    - results of, 137
  - SQLFire, comparing with RDBMS, 121-125
  - vFabric Reference Architecture, 210-211
  - vSphere 5, enhancements to, 142-143
- policies (GC), selecting, 184
- IBM JVM, 186
- Oracle HotSpot JVM, 184-185
- Oracle jRockit JVM, 187
- pools, best practices, 158-160
- practical sizing limits of Java platforms, 4-6
- production environment, sizing, 95
- properties of load profile, 80-81

## Q-R

- queries, entity groups, 97
- RAM
  - available RAM per NUMA node, calculating, 10-12
  - theoretical and practical sizing limits of Java, 5
- RDBMS performance, comparing to SQLFire, 121-123
  - CPU utilization, 125
  - response time, 124
  - scalability, 124
- redundancy (SQLFire), 31, 37-38
- redundancy zones (SQLFire), 28
- resource pools, best practices, 158-160
- response time of SQLFire, comparing with RDBMS, 124
- resxtp vCLI reference, 198

- Rockit, best practices, 181-183
- Rosenberg, Harold, 127
- RowLoader (SQLFire), 46-47

## S

- scalability
  - benchmarking environment, sizing, 91, 95
  - horizontal scalability
    - best practices, 158-160
    - testing, 86-91
    - vSphere features, 218
  - JVMs and VMs, 16-20
  - production environment, sizing, 95
  - SQLFire
    - comparing with RDBMS, 124
    - sizing example, 112, 115-119
  - vertical scalability
    - best practices, 156-157
    - testing, 82-83, 86
    - vSphere features, 218
- ScavengeBeforeFullGC JVM configuration option, 67
- selecting
  - daily transactional data, 97
  - GC policies, 184
    - IBM JVM, 186
    - Oracle HotSpot JVM, 184-185
    - Oracle jRockit, 187
- server groups (SQLFire), 30-34
- server tier, SQLFire client/server topology, 25
- shared-nothing persistence mechanism (SQLFire), 22
- sizing
  - benchmarking environment, 91, 95
  - category 1 workloads, VM sizing and configuration, 148
  - JVM, 227-228
  - large VMs and JVMs, NUMA implications, 108-111
  - production environment, 95
  - SQLFire Java platforms
    - category 2 workloads, 97-100, 104-106
    - sizing example, 112, 115-119
- SMA (Spring Migration Analyzer), 205-207
- Spring Framework, 214
- Spring Integration, 215
- Spring Tools Suite, 213
- Spring Travel
  - comparing SQLFire performance versus RDBMS, 121-125
  - downloading, 122
- SpringTrader performance study, 131-132
  - application and data tier, 133-137
  - results of, 137

- SQL queries, entity groups, 97
- SQLFire, 22, 209-211, 216
  - active-active architectures, 57, 60
  - Apache Derby project, 23
  - asynchronous listeners, 52-54
  - best practices, 166-172
  - client/server topology, 24
    - best practices, 26
    - client tier, 25
    - server tier, 25
  - colocation, 38-39
  - DBSynchronizer, 54-56
  - DDLUtils, 57
  - disk persistence, 39-41
  - enterprise data fabric system, 24
  - features, 22, 32
  - global multisite topology, 28-30
  - listeners, 47-49
  - low-end volume, 24
  - number of members, calculating, 105-106
  - on vSphere best practices, 173-181
  - partitioning, 31, 34-37
  - peer-to-peer topology, 27
  - performance, comparing with RDBMS, 121-125
  - redundancy, 31, 37-38
  - redundancy zones, 28
  - RowLoader, 46-47
  - server groups, 30-34
  - shared-nothing persistence mechanism, 22
  - sizing Java platforms, category 2 workloads, 97-100, 104-106
  - transactions, 41-46
  - writers, 50-51
- stacking, JVM stacking, 83
- statistics, viewing performance in vCenter, 195
- support for IBM WebSphere Application Server, 221
- support for VMware, 209
- support for WebLogic, 221
- support-request tickets, opening, 191
- survivor spaces tuning (GC), 78
- SurvivorRatio JVM configuration option, 67
- SYN cookies, troubleshooting, 179

## T

- TargetSurvivorRatio=80 JVM configuration option, 67
- technical considerations, large-scale Java platforms, 4
- technology tiers, 189-190
- theoretical sizing limits of Java platforms, 4-6
- threading
  - HT, 83, 224

- thread contentions, troubleshooting Java, 203-204
- thread dumps, inspecting, 204
- thread ratios, establishing, 93
- throughput GC, 63-64
- tiers
  - inter-tier configuration, best practices, 160-164
  - migration phases, 220
  - virtualizing, 220
- timekeeping, VM timekeeping best practices, 156
- topologies (SQLFire)
  - client/server topology, 24
    - best practices, 26
    - client tier, 25
    - server tier, 25
  - global multisite topology, 28-30
  - peer-to-peer topology, 27
  - redundancy zones, 28
- transactions (SQLFire), 41-46
- troubleshooting
  - collecting metrics from vCenter, 191-193
  - CPUvCenter performance charts, 193
  - esxtop, 195, 198
  - Java, 198, 201
    - memory, 202-203
    - thread contentions, 203-204
  - performance charts, accessing, 195
  - support-request tickets, opening, 191
  - vCenter
    - disk usage charts, 194
    - network performance charts, 194
- tuning
  - GC
    - old generation tuning, 76
    - survivor spaces tuning, 78
    - young generation tuning, 71-74
  - knowledge requirements for large-scale Java platforms, 7-9
  - monitoring system, 225-226

## U

- UseBiasedLocking JVM configuration option, 67
- UseCMSInitiatingOccupancyOnly JVM configuration option, 67
- UseCompressedOops JVM configuration option, 69
- UseCompressedStrings JVM configuration option, 69
- UseConcMarkSweepGC JVM configuration option, 66
- UseNUMA JVM configuration option, 69
- UseParNewGC JVM configuration option, 67
- UseStringCache JVM configuration option, 69

## V

### vCenter

- CPU performance charts, 193
- disk usage charts, 194
- memory counters, 192
- metrics, collecting, 191-193
- network performance charts, 194

### vCPU for VM best practices, 149-150

### vertical scalability

- best practices, 156-157
- JVMs and VMs, 16-20
- tests, establishing benchmark for, 82-83, 86
- vSphere features, 218

### vFabric

- Application Director, 217
- Data Director, 217
- RabbitMQ, 216
- Spring Framework, 214
- Spring Integration, 215
- Spring Tools Suite, 213
- SQLFire. *See* SQLFire
- tc Server, 215
- vPostgres, 216
- Web Server, 215

### vFabric EM4J (Elastic Memory for Java), 215

### vFabric Reference Architecture, 110, 210-211

### viewing performance statistics, 195

### virtualizing WebLogic application servers, 212

### VisualVM, 201

### vMA (vSphere Management Assistant), 198

### VMotion, 142

### VMs (virtual machines), 3

- building block VM, 83
- category 1 workloads
  - high-level vSphere best practices, 165-166
  - horizontal scalability best practices, 158-160
  - inter-tier configuration, 160-164
  - sizing and configuration, 148
  - vCPU for VM best practices, 149-150
  - vertical scalability best practices, 156-157
  - VM memory size best practices, 150-155
  - VM timekeeping best practices, 156
- category 2 Java workloads, 27, 145
  - best practices, 166-181
  - sizing SQLFire Java platforms, 97-100, 104-106
- category 3 workloads, 145, 181-183
  - client/server topology (SQLFire), 25
  - memory requirements, calculating, 84-86
  - number of JVMs for, calculating, 223
  - scaling, 223
  - sizing, NUMA implications, 108-111

### VMware support, 209

### VMware vCenter Site Recovery Manager, 220

### vNUMA, 15, 110

### vSphere

- application servers, recommendations for, 219
- high-level vSphere, best practices, 165-166
- Java applications
  - horizontal scalability, 218
  - performance, 207
  - vertical scalability, 218
- Java coding practices, 217
- licensing, 221
- performance, 209
- troubleshooting with esxtop, 195, 198
- WebLogic, running, 219

### vSphere 5, performance enhancements, 142-143

## W

### web server tier, 190

### WebLogic, 219

- licensing, 221
- support for, 221
- VMware customer references, 222

### workloads

#### category 1, 145

- high-level vSphere best practices, 165-166
- horizontal scalability best practices, 158-160
- inter-tier configuration best practices, 160-164
- vCPU for VM best practices, 149-150
- vertical scalability best practices, 156-157
- VM memory size best practices, 150-155
- VM sizing and configuration best practices, 148

#### VM timekeeping best practices, 156

#### category 2, 27, 145, 166-167

- sizing SQLFire Java platforms, 97-100, 104-106

#### SQLFire best practices, 166-172

#### SQLFire on vSphere best practices, 173-181

#### category 3, 145

- client/server topology (SQLFire), 25
- IBM JVM and Oracle Rockit JVM, 181-183
- Olio performance study, 127, 131
- profile, establishing, 80-81
- properties, 80-81
- SpringTrader performance study, 131-137

### writers (SQLFire), 32, 50-51

## X-Y-Z

### -Xmn, adjusting young generation internal cycle, 72-74

### -Xmn21g JVM configuration option, 66

### young generation tuning (GC), 71-74