

ADDISON
WESLEY
DATA &
ANALYTICS
SERIES



PRACTICAL Cassandra

A Developer's Approach

RUSSELL BRADBERRY
ERIC LUBOW

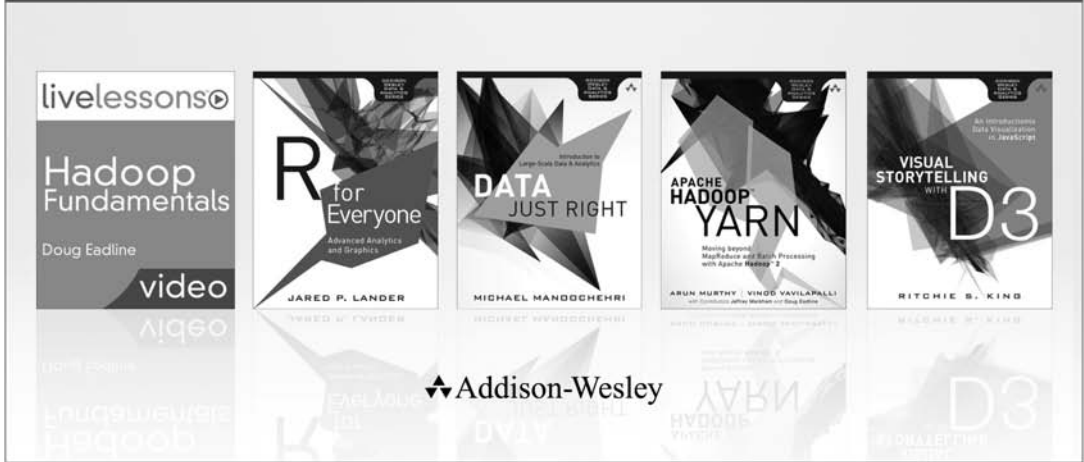
FREE SAMPLE CHAPTER



SHARE WITH OTHERS

Practical Cassandra

The Addison-Wesley Data and Analytics Series



Visit informit.com/awdataseries for a complete list of available publications.

The Addison-Wesley Data and Analytics Series provides readers with practical knowledge for solving problems and answering questions with data. Titles in this series primarily focus on three areas:

1. **Infrastructure:** how to store, move, and manage data
2. **Algorithms:** how to mine intelligence or make predictions based on data
3. **Visualizations:** how to represent data and insights in a meaningful and compelling way

The series aims to tie all three of these areas together to help the reader build end-to-end systems for fighting spam; making recommendations; building personalization; detecting trends, patterns, or problems; and gaining insight from the data exhaust of systems and user interactions.



Make sure to connect with us!
informit.com/socialconnect

informit.com
the trusted technology learning source

Addison-Wesley

Safari
Books Online

Practical Cassandra

A Developer's Approach

Russell Bradberry
Eric Lubow

◆◆Addison-Wesley

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

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

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

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

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

For questions about sales outside the U.S., please contact international@pearsoned.com.

Visit us on the Web: informit.com/aw

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

Copyright © 2014 Pearson Education, Inc.

All rights reserved. Printed in the United States of America. This publication is protected by copyright, and permission must be obtained from the publisher prior to any prohibited reproduction, storage in a retrieval system, or transmission in any form or by any means, electronic, mechanical, photocopying, recording, or likewise. To obtain permission to use material from this work, please submit a written request to Pearson Education, Inc., Permissions Department, One Lake Street, Upper Saddle River, New Jersey 07458, or you may fax your request to (201) 236-3290.

ISBN-13: 978-0-321-93394-2

ISBN-10: 0-321-93394-X

Text printed in the United States on recycled paper at RR Donnelley in Crawfordsville, Indiana. First printing, December 2013



This book is for the community. We have been a part of the Cassandra community for a few years now, and they have been fantastic every step of the way. This book is our way of giving back to the people who have helped us and have allowed us to help pave the way for the future of Cassandra.



This page intentionally left blank

Contents

Foreword by Jonathon Ellis xiii

Foreword by Paul Dix xv

Preface xvii

Acknowledgments xxi

About the Authors xxiii

1 Introduction to Cassandra 1

A Greek Story 1

What Is NoSQL? 2

There's No Such Thing as "Web Scale" 2

ACID, CAP, and BASE 2

ACID 3

CAP 3

BASE 4

Where Cassandra Fits In 5

What Is Cassandra? 5

History of Cassandra 6

Schema-less (If You Want) 7

Who Uses Cassandra? 7

Is Cassandra Right for Me? 7

Cassandra Terminology 8

Cluster 8

Homogeneous Environment 8

Node 8

Replication Factor 8

Tunable Consistency 8

Our Hope 9

2 Installation 11

Prerequisites 11

Installation 11

Debian 12

RedHat/CentOS/Oracle 12

From Binaries 12

Configuration	13
Cluster Setup	15
Summary	16

3 Data Modeling 17

The Cassandra Data Model	17
Model Queries—Not Data	19
Collections	22
Sets	22
Lists	23
Maps	24
Summary	25

4 CQL 27

A Familiar Way of Doing Things	27
CQL 1	27
CQL 2	28
CQL 3	28
Data Types	28
Commands	30
Example Schemas	37
Summary	39

5 Deployment and Provisioning 41

Keyspace Creation	41
Replication Factor	41
Replication Strategies	42
SimpleStrategy	42
NetworkTopologyStrategy	42
Snitches	43
Simple	43
Dynamic	43
Rack Inferring	44
EC2	44
Ec2MultiRegion	45
Property File	45
PropertyFileSnitch Configuration	46
Partitioners	46

Byte Ordered	47
Random Partitioners	47
Node Layout	48
Virtual Nodes	48
Balanced Clusters	49
Firewalls	49
Platforms	49
Amazon Web Services	50
Other Platforms	50
Summary	50

6 Performance Tuning 51

Methodology	51
Testing in Production	52
Tuning	52
Timeouts	52
CommitLog	53
MemTables	54
Concurrency	55
Durability and Consistency	55
Compression	56
SnappyCompressor	58
DeflateCompressor	58
File System	58
Caching	59
How Cassandra Caching Works	59
General Caching Tips	59
Global Cache Tuning	60
ColumnFamily Cache Tuning	61
Bloom Filters	61
System Tuning	62
Testing I/O Concurrency	62
Virtual Memory and Swap	63
sysctl Network Settings	64
File Limit Settings	64
Solid-State Drives	64
JVM Tuning	65

Multiple JVM Options	65
Maximum Heap Size	65
Garbage Collection	66
Summary	67

7 Maintenance 69

Understanding <code>nodetool</code>	69
General Usage	71
Node Information	72
Ring Information	72
ColumnFamily Statistics	73
Thread Pool Statistics	74
Flushing and Draining	75
Cleaning	75
<code>upgradesstables</code> and <code>scrub</code>	76
Compactions	76
What, Where, Why, and How	76
Compaction Strategies	77
Impact	77
Backup and Restore	79
Are Backups Necessary?	79
Snapshots	79
CommitLog Archiving	81
<code>archive_command</code>	81
<code>restore_command</code>	81
<code>restore_directories</code>	81
<code>restore_point_in_time</code>	82
CommitLog Archiving Notes	82
Summary	82

8 Monitoring 83

Logging	83
Changing Log Levels	84
Example Error	84
JMX and MBeans	85
JConsole	86
Health Checks	91

Nagios	91
Cassandra-Specific Health Checks	94
Cassandra Interactions	96
Summary	96

9 Drivers and Sample Code 99

Java	100
C#	104
Python	108
Ruby	112
Summary	117

10 Troubleshooting 119

Toolkit	119
iostat	119
dstat	120
nodetool	121
Common Problems	121
Slow Reads, Fast Writes	122
Freezing Nodes	123
Tracking Down OOM Errors	124
Ring View Differs between Nodes	124
Insufficient User Resources	124
Summary	126

11 Architecture 127

Meta Keyspaces	127
System Keyspace	127
Authentication	128
Gossip Protocol	129
Failure Detection	130
CommitLogs and MemTables	130
SSTables	130
HintedHandoffs	131
Bloom Filters	131
Compaction Types	132
Tombstones	132
Staged Event-Driven Architecture	133
Summary	134

12 Case Studies 135

Ooyala 135

Hailo 137

Taking the Leap 138

Proof Is in the Pudding 139

Lessons Learned 140

Summary 141

eBay 141

eBay's Transactional Data Platform 141

Why Cassandra? 142

Cassandra Growth 143

Many Use Cases 143

Cassandra Deployment 146

Challenges Faced and Lessons Learned 147

Summary 147

A Getting Help 149

Preparing Information 149

IRC 149

Mailing Lists 149

B Enterprise Cassandra 151

DataStax 151

Acunu 152

Titan by Aurelius 153

Pentaho 154

Instaclustr 154

Index 157

Foreword by Jonathon Ellis

I was excited to learn that *Practical Cassandra* would be released right at my five-year anniversary of working on Cassandra. During that time, Cassandra has achieved its goal of offering the world's most reliable and performant scalable database. Along the way, Cassandra has changed significantly, and a modern book is, at this point, overdue. Eric and Russell were early adopters of Cassandra at SimpleReach; in *Practical Cassandra*, you benefit from their experience in the trenches administering Cassandra, developing against it, and building one of the first CQL drivers.

If you are deploying Cassandra soon, or you inherited a Cassandra cluster to tend, spend some time with the deployment, performance tuning, and maintenance chapters. Some complexity is inherent in a distributed system, particularly one designed to push performance limits and scale without compromise; forewarned is, as they say, forearmed. If you are new to Cassandra, I highly recommend the chapters on data modeling and CQL. The Cassandra Query Language represents a major shift in developing against Cassandra and dramatically lowers the learning curve from what you may expect or fear.

Here's to the next five years of progress!

—Jonathon Ellis, *Apache Cassandra Chair*

This page intentionally left blank

Foreword by Paul Dix

Cassandra is quickly becoming one of the backbone components for anyone working with large datasets and real-time analytics. Its ability to scale horizontally to handle hundreds of thousands (or millions) of writes per second makes it a great choice for high-volume systems that must also be highly available. That's why I'm very pleased that this book is the first in the series to cover a key infrastructural component for the Addison-Wesley Data & Analytics Series: the data storage layer.

In 2011, I was making my second foray into working with Cassandra to create a high-volume, scalable time series data store. At the time, Cassandra 0.8 had been released, and the path to 1.0 was fairly clear, but the available literature was lagging sorely behind. This book is exactly what I could have used at the time. It provides a great introduction to setting up and modeling your data in Cassandra. It has coverage of the most recent features, including CQL, sets, maps, and lists. However, it doesn't stop with the introductory stuff. There's great material on how to run a cluster in production, how to tune performance, and on general operational concerns.

I can't think of more qualified users of Cassandra to bring this material to you. Eric and Russell are Datastax Cassandra MVPs and have been working extensively with Cassandra and running it in production for years. Thankfully, they've done a great job of distilling their experience into this book so you won't have to search for insight into how to develop against and run the most current release of Cassandra.

—Paul Dix, *Series Editor*

This page intentionally left blank

Preface

Apache Cassandra is a massively scalable, open-source, NoSQL database. Cassandra is best suited to applications that need to store large amounts of structured, semistructured, and unstructured data. Cassandra offers asynchronous masterless replication to nodes in many data centers. This gives it the capability to have no single point of failure while still offering low latency operations.

When we first embarked on the journey of writing a book, we had one goal in mind: We wanted to keep the book easily digestible by someone just getting started with Cassandra, but also make it a useful reference guide for day-to-day maintenance, tuning, and troubleshooting. We know the pain of scouring the Internet only to find outdated and contrived examples of how to get started with a new technology. We hope that *Practical Cassandra* will be the go-to guide for developers—both new and at an intermediate level—to get up and running with as little friction as possible.

This book describes, in detail, how to go from nothing to a fully functional Cassandra cluster. It shows how to bring up a cluster of Cassandra servers, choose the appropriate configuration options for the cluster, model your data, and monitor and troubleshoot any issues. Toward the end of the book, we provide sample code, in-depth detail as to how Cassandra works under the covers, and real-world case studies from prominent users.

What's in This Book?

This book is intended to guide a developer in getting started with Cassandra, from installation to common maintenance tasks to writing an application. If you are just starting with Cassandra, this book will be most helpful when read from start to finish. If you are familiar with Cassandra, you can skip around the chapters to easily find what you need.

- **Chapter 1, Introduction to Cassandra:** This chapter gives an introduction to Cassandra and the philosophies and history of the project. It provides an overview of terminology, what Cassandra is best suited for, and, most important what we hope to accomplish with this book.
- **Chapter 2, Installation:** Chapter 2 is the start-to-finish guide to getting Cassandra up and running. Whether the installation is on a single node or a large cluster, this chapter guides you through the process. In addition to cluster setup, the most important configuration options are outlined.
- **Chapter 3, Data Modeling:** Data modeling is one of the most important aspects of using Cassandra. Chapter 3 discusses the primary differences between Cassandra

and traditional RDBMSs, as well as going in depth into different design patterns, philosophies, and special features that make Cassandra the data store of tomorrow.

- Chapter 4, CQL: CQL is Cassandra’s answer to SQL. While not a full implementation of SQL, CQL helps to bridge the gap when transitioning from an RDBMS. This chapter explores in depth the features of CQL and provides several real-world examples of how to use it.
- Chapter 5, Deployment and Provisioning: After you’ve gotten an overview of installation and querying, this chapter guides you through real-world deployment and resource provisioning. Whether you plan on deploying to the cloud or on bare-metal hardware, this chapter is for you. In addition to outlining provisioning in various types of configurations, it discusses the impact of the different configuration options and what is best for different types of workloads.
- Chapter 6, Performance Tuning: Now that you have a live production cluster deployed, this chapter guides you through tweaking the Cassandra dials to get the most out of your hardware, operating system, and the Java Virtual Machine (JVM).
- Chapter 7, Maintenance: Just as with everything in life, the key to having a performant and, more important, working Cassandra cluster is to maintain it properly. Chapter 7 describes all the different tools that take the headache out of maintaining the components of your system.
- Chapter 8, Monitoring: Any systems administrator will tell you that a healthy system is a monitored system. Chapter 8 outlines the different types of monitoring options, tools, and what to look out for when administering a Cassandra cluster.
- Chapter 9, Drivers and Sample Code: Now that you have a firm grasp on how to manage and maintain your Cassandra cluster, it is time to get your feet wet. In Chapter 9, we discuss the different drivers and driver features offered in various languages. We then go for the deep dive by presenting a working example application in not only one, but four of the most commonly used languages: Java, C#, Ruby, and Python.
- Chapter 10, Troubleshooting: Now that you have written your sample application, what happens when something doesn’t quite work right? Chapter 10 outlines the tools and techniques that can be used to get your application back on the fast track.
- Chapter 11, Architecture: Ever wonder what goes on under the Cassandra “hood”? In this chapter, we discuss how Cassandra works, how it keeps your data safe and accurate, and how it achieves such blazingly fast performance.
- Chapter 12, Case Studies: So who uses Cassandra, and how? Chapter 12 presents three case studies from forward-thinking companies that use Cassandra in unique ways. You will get the perspective straight from the mouths of the developers at Ooyala, Hailo, and eBay.

- Appendix A, Getting Help: Whether you're stuck on a confusing problem or just have a theoretical question, having a place to go for help is paramount. This appendix tells you about the best places to get that help.
- Appendix B, Enterprise Cassandra: There are many reasons to use Cassandra, but sometimes it may be better for you to focus on your organization's core competencies. This appendix describes a few companies that can help you leverage Cassandra efficiently and effectively while letting you focus on what you do best.

Code Samples

All code samples and more in-depth examples can be found on GitHub at <http://devdazed.github.io/practical-cassandra/>.

This page intentionally left blank

Acknowledgments

We would like to acknowledge everyone involved with Cassandra and the Cassandra community—everyone from the core contributors of Cassandra all the way down to the end users who have made it such a popular platform to work with. Without the community, Cassandra wouldn't be where it is today. Special thanks go to

- Jay Patel for putting together the eBay case study
- Al Tobey and Evan Chan for putting together the case study on Ooyala
- Dominic Wong for putting together the Hailo case study
- All the technical reviewers, including Adam Chalemian, Mark Herschberg, Joe Stein, and Bryan Smith, who helped give excellent feedback and ensured technical accuracy where possible
- Paul Dix for setting us up and getting us on the right track with writing

This page intentionally left blank

About the Authors

Russell Bradberry (Twitter: @devdazed) is the principal architect at SimpleReach, where he is responsible for designing and building out highly scalable, high-volume, distributed data solutions. He has brought to market a wide range of products, including a real-time bidding ad server, a rich media ad management tool, a content recommendation system, and, most recently, a real-time social intelligence platform. He is a U.S. Navy veteran, a DataStax MVP for Apache Cassandra, and the author of the NodeJS Cassandra driver Helenus.

Eric Lubow (Twitter: @elubow) is currently chief technology officer of SimpleReach, where he builds highly scalable, distributed systems for processing social data. He began his career building secure Linux systems. Since then he has worked on building and administering various types of ad systems, maintaining and deploying large-scale Web applications, and building email delivery and analytics systems. He is also a U.S. Army combat veteran and a DataStax MVP for Apache Cassandra.

Eric and Russ are regular speakers about Cassandra and distributed systems, and both live in New York City.

This page intentionally left blank

Monitoring

As the old systems adage goes, a service doesn't exist unless it's monitored. In this chapter, we will cover the basics of monitoring Cassandra. These include file-based logging, inspection of the JVM, and monitoring of Cassandra itself.

Logging

Under the covers, Cassandra uses the standard Java log library Log4j. Log4j is another Apache project that enables the capability to control the granularity of log statements using a configuration file. If you want to find out more about what is happening on a particular node than what `nodetool` and JMX MBeans (which we will cover in more detail later in the chapter) are telling you, you can change the logging levels.

As a front end to the Log4j back end, Cassandra uses Simple Logging Façade for Java (SLF4J). The logging levels from least verbose to most verbose are

- TRACE
- DEBUG
- INFO
- WARN
- ERROR
- FATAL

Understanding these logging levels is important not only to help monitor what is going on in the system or on a particular node but also to help troubleshoot problems. In troubleshooting complex systems such as Cassandra, Cassandra's `nodetool`, logging, and even the JMX MBeans can lead to red herrings. So it is necessary to compile as much information pertinent to the problem as possible to help diagnose what might be going on.

Taking a look at a normal healthy Cassandra node's `system.log`, you will see `INFO` lines that refer to various stages of the system executing their tasks. These include MemTable flushes, HintedHandoffs, and compactions, just to name a few.

Changing Log Levels

If you want to make any changes to the logging schema, you will need to find the `log4j-server.properties` file. The default logging level for Cassandra and the `rootLogger` is `INFO`. This level provides a standard amount of information that is sufficient for understanding the general health of your system. It is definitely helpful to see what your system looks like, so you should do so while logging at the `DEBUG` level. Be sure not to leave Cassandra in `DEBUG` mode for production as the entire system will act noticeably slower. To change the standard logging level in Cassandra from `INFO` to `DEBUG`, change the line that looks like this:

```
log4j.rootLogger=INFO,stdout,R
```

to this:

```
log4j.rootLogger=DEBUG,stdout,R
```

Now your Cassandra node will be running in `DEBUG` mode. To change it back, just swap the `INFO` and `DEBUG` again. To show less logging, you can change the logging level to `WARN`, `ERROR`, or `FATAL`.

Example Error

It is worth noting that not all problem messages enter the logs as `WARNING` or higher (higher meaning toward `FATAL`). Listing 8.1 presents an example of when things start to go south. This is a common set of log messages that you may see with your system set at `INFO` level. Even with the logging level set to `INFO`, there is a lot of useful information in the logs. Don't be afraid to keep a regular eye on the logs so you know what patterns of log messages are normal for your system. For example, if things are starting to slow down, you may see something like Listing 8.1.

Listing 8.1 INFO Messages That Show Mutation and READ Messages Dropped

```
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,290 MessagingService.java (line 607)
3476 MUTATION messages dropped in last5000ms
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,290 MessagingService.java (line 607)
677 READ messages dropped in last 5000ms
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,291 StatusLogger.java (line 50)
Pool Name           Active           Pending          Blocked
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,291 StatusLogger.java (line 65)
ReadStage           32              621             0
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,291 StatusLogger.java (line 65)
RequestResponseStage 0                0               0
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,291 StatusLogger.java (line 65)
ReadRepairStage     0                0               0
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,291 StatusLogger.java (line 65)
MutationStage       32              4105            0
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,291 StatusLogger.java (line 65)
ReplicateOnWriteStage 0                0               0
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,292 StatusLogger.java (line 65)
GossipStage         0                0               0
```

```

INFO [ScheduledTasks:1] 2012-07-09 20:48:57,292 StatusLogger.java (line 65)
AntiEntropyStage           0           0           0
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,292 StatusLogger.java (line 65)
MigrationStage             0           0           0
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,292 StatusLogger.java (line 65)
StreamStage                0           0           0
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,292 StatusLogger.java (line 65)
MemtablePostFlusher       0           0           0
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,292 StatusLogger.java (line 65)
FlushWriter                0           0           0
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,292 StatusLogger.java (line 65)
MiscStage                  0           0           0
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,292 StatusLogger.java (line 65)
InternalResponseStage     0           0           0
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,293 StatusLogger.java (line 65)
HintedHandoff              1           8           0
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,293 StatusLogger.java (line 70)
CompactionManager          0           0
INFO [ScheduledTasks:1] 2012-07-09 20:48:57,293 StatusLogger.java (line 82)
MessagingService           n/a         0,0

```

When you see messages being dropped, as in the first two lines of the listing, that's a sign that your system is under stress. Depending on your application requirements, some level of dropped messages may be acceptable. But regardless of whether or not your application can tolerate running in a degraded state, the overall health of your cluster (or at the very least this node) is in question. Most applications are capable of handling dropped READ requests, but dropped MUTATION messages mean that data that should have been written isn't getting written. Depending on the consistency level of the write in question, this could mean the write didn't happen at all, or it could mean the write didn't happen on this node. Also notice that the `ReadStage` and `MutationStage` lines have multiple `Active` and `Pending` messages left to work on. The reason these messages are dropped is that Cassandra wants to do its best to keep up with the volume of work that it is being given.

There are other such common log lines to watch for, which can be done via a log monitor. One method for monitoring the logs programmatically using Nagios will be discussed later in this chapter.

JMX and MBeans

Built into Cassandra and the JVM is the capability to use the JMX, or Java Management Extensions. In other words, using JMX gives you the capability to manage your servers remotely or check into settings programmatically, including the memory, CPU, threads, Gossip, or any other part of the system that has been instrumented in JMX. Instrumentation is what enables the application to provide application-specific

information to be collected by external tools. JMX also gives you the ability to control certain aspects of this information.

JConsole

MBeans, or Managed Beans, are a special type of JavaBean that makes a resource inside the application or the JVM available externally. The standard tool that ships with Java for managing the MBeans is JConsole.

The most common use case for accessing a Cassandra server is that the server will be remote and you probably won't have console access to it. It is also highly recommended that you run JConsole remotely as it is a heavy user of resources on a machine and can steal those resources away from the Cassandra node. If this is the case, you can use SSH tunneling to bring up JConsole. When you SSH, be sure to use the `-X` switch to ensure that X11 forwarding is on. This is what enables you to use JConsole over the network. After SSHing into the machine running Cassandra and finding the JConsole binary, just execute it as you would any normal binary. Assuming everything is configured correctly, you will get a JConsole login window as shown in Figure 8.1.

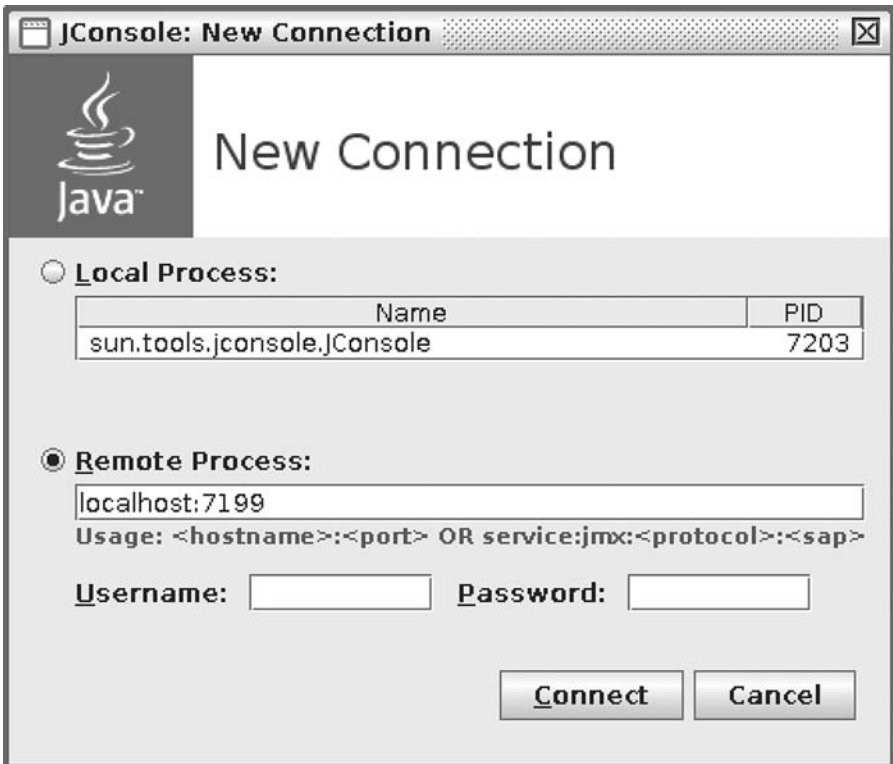


Figure 8.1 JConsole login window when logging in via localhost

Click on the radio button labeled Remote Process and type `localhost:7199`. If you have Cassandra set up with authentication, you will need to put in the username and password as well. Port 7199 is the default JMX port for Cassandra. The first thing you will notice once a connection has been established is that there are multiple tabs that contain information for you to look through. These tabs are Overview, Memory, Threads, Classes, VM Summary, and MBeans.

The Overview, Memory, and Threads tabs are sets of graphs that provide insight into the current state of the system. The Classes graph is simply a graph of how many classes are loaded into the JVM at the current time (or a different time range that you choose). The VM Summary is an overview of what the current view of the Java Virtual Machine is.

Memory

The Memory tab consists of graphs about the current state of the memory (see Figure 8.2). One of the most important memory stats that you want to be aware of is your current heap usage. It is also the default graph that comes up on the Memory tab. As with everything else in your system, it is helpful to know what a good baseline is for heap usage during normal system operations. If you have a 10GB heap set and you find during most of your operations that you are using only 3GB, you can likely reduce your JVM maximum heap size. This will enable you to reduce your memory footprint on the system and possibly even speed up your GCs.

Cassandra also does a lot of work off-heap. Things like bloom filters in version 1.2 and forward are off-heap now. Keeping tabs on what the off-heap memory usage looks like is also very important.

Garbage collection is also one of the metrics that can be viewed from the Memory tab. If needed (though typically not recommended unless you know what you are doing), you can even force a GC from the Memory tab in JConsole. In Cassandra 1.1 and later, there are even helpful bars that display the total memory used versus memory available both on-heap and off-heap.

Threads

The Threads tab in JConsole is dedicated to showing the current and peak usage patterns of various thread stages in Cassandra (see Figure 8.3). These include everything that you would normally see in the logs from things like the CommitLog handler and compaction scheduler all the way to garbage collection, Gossip, and streaming. It is also helpful here to see how many threads your system uses normally as well as under load.

MBeans

The final tab in JConsole is the tab for MBeans (see Figure 8.4). There are a lot of MBeans that are useful for assessing the state of your Cassandra node and your Cassandra cluster. You will notice that there are a few groupings here that can be expanded. Other than the standard Java MBeans that are available to every agent, there are several groupings specific to Cassandra. All of their class paths start with `org.apache.cassandra.`

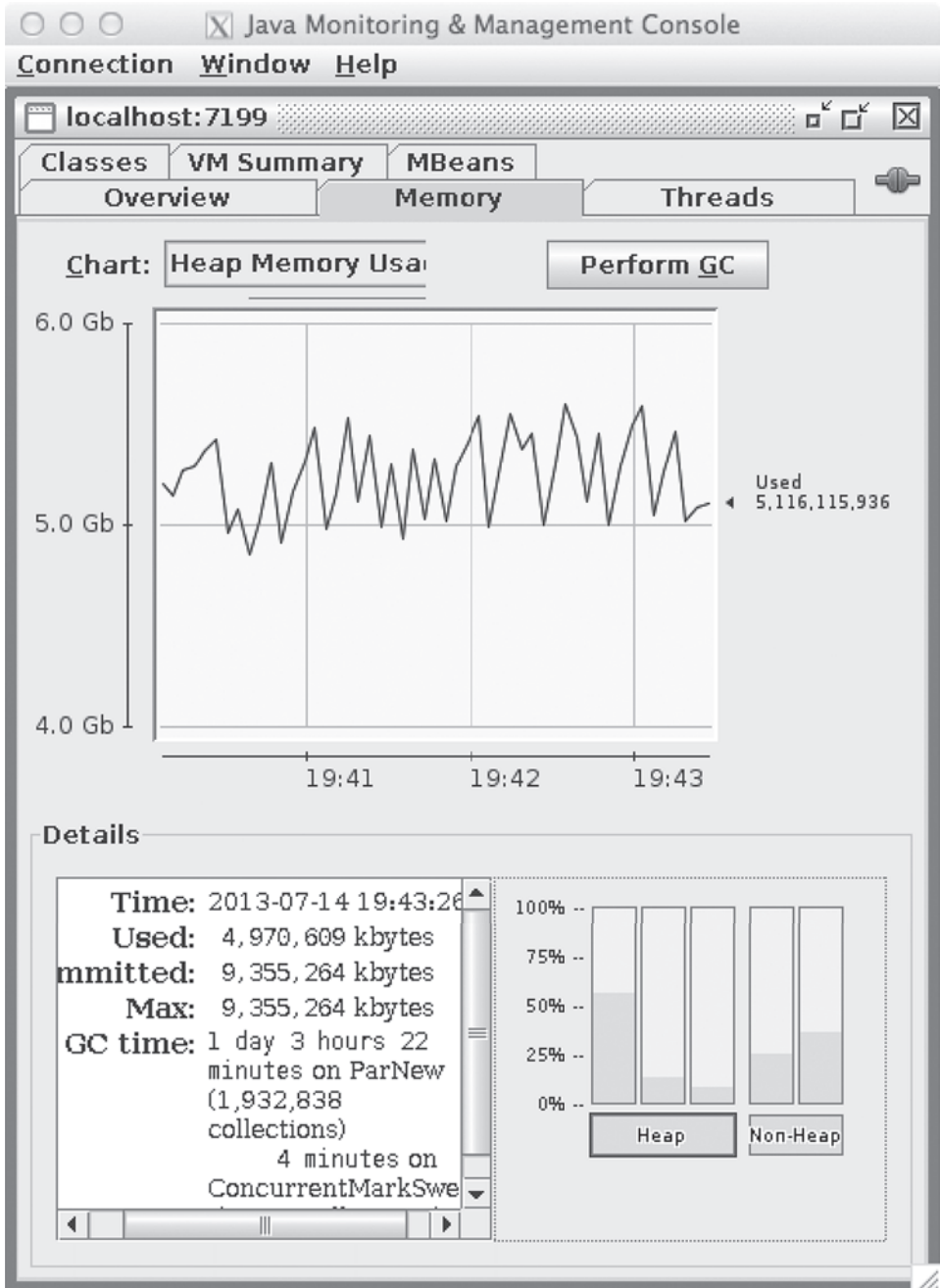


Figure 8.2 JConsole Memory tab displaying heap usage graphs

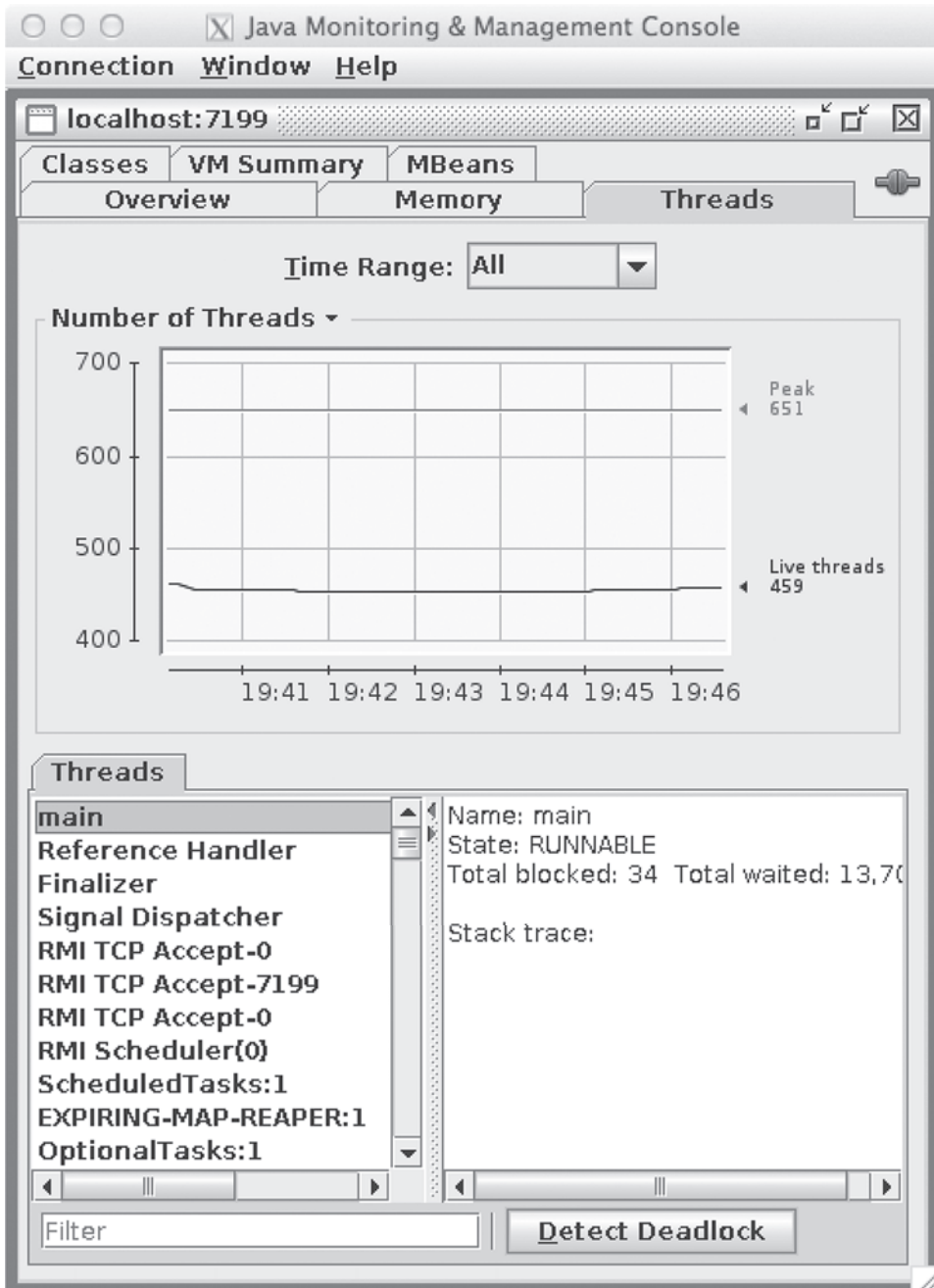


Figure 8.3 JConsole Threads tab displaying the number of threads Cassandra is currently using and general information on the main Cassandra thread

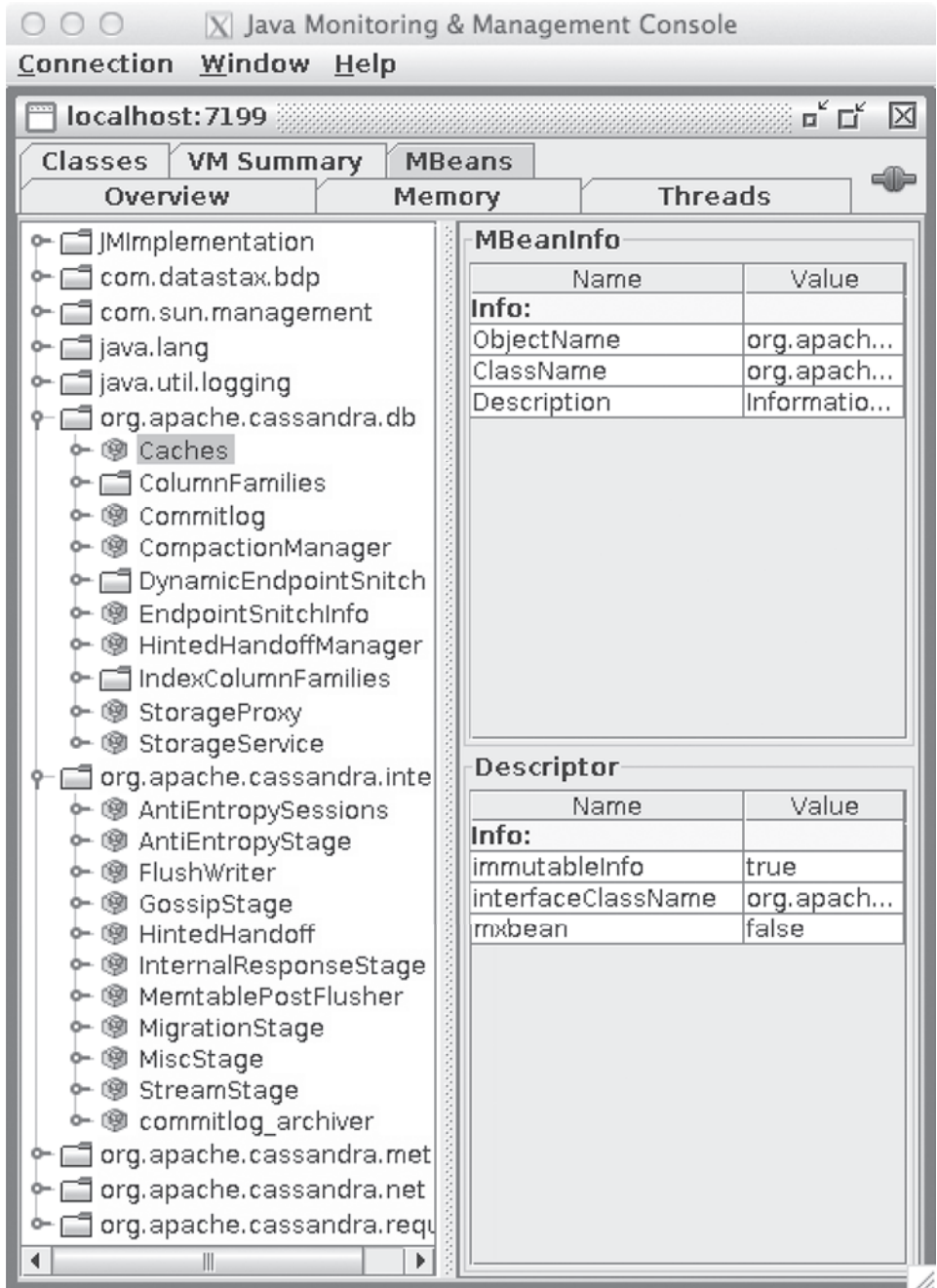


Figure 8.4 JConsole MBeans tab view with the Cassandra DB and Internal trees opened

There are a lot of MBeans provided by an application as complex as Cassandra. There is no need to cover all of them as you can easily explore them on your own using the JConsole interface.

As a high-level overview, they are broken down into the following categories:

- **DB.** The MBeans stored in the DB section cover everything about the actual data storage part of Cassandra. You can view information about the cache and CommitLogs, or even information about the individual ColumnFamilies that you have created in each of the keyspaces. HintedHandoffManager, EndPointSnitchInfo, and CompactionManager information can also be found here.
- **Internal.** In the Internal section, there are MBeans that cover the state and statistics around the staged architecture. More specifically, you can find information about the state of Gossip and HintedHandoff as opposed to just finding information about the managers as in the DB section.
- **Metrics.** The metrics available in this section are ClientRequestMetrics. These are things like read and write timeouts and Unavailable errors.
- **Net.** The network section houses information about internode communication. This includes information about the FailureDetector, Gossiper, MessagingService, and StreamingService.
- **Request.** The MBeans in the Request section are about tasks related to read, write, and replication.

When you select any MBean in the tree, its MBeanInfo and MBean descriptor are displayed on the right-hand side of the window. If any additional attributes, operations, or notifications are available, they will appear in the tree as well, below the selected MBean.

Each one of these sections of MBeans provides access to a large amount of information, giving you insight into both the system as a whole and the individual nodes. Familiarizing yourself with what is available to you as an admin will help you when it comes time to instrument JMX-level checks from within your monitoring system.

Health Checks

Using JConsole to monitor your system is tedious and good as a monitoring system only if you are actively staring at the graphs and information all the time. Since that is unrealistic and time-consuming, we recommend that you use other systems for monitoring the general health of your system such as Nagios.

Nagios

Nagios is open-source software dedicated to monitoring computers, networks, hosts, and services and can alert you when things are going wrong or have been resolved. It is extremely versatile and has the capability to monitor many types of services, applications, or parts of an application. Let's start at the bottom of the monitoring chain and work our way up. In order to avoid a complete lesson on monitoring, we will only cover the basics along with what the most common checks should be as they relate to Cassandra and its operation.

There are three primary alerts in Nagios: `WARNING`, `CRITICAL`, and `OK`. They mean exactly what they sound like. A `WARNING` alert is sent if the service in question is starting to show signs of a problem, such as a hard drive nearing capacity. A `CRITICAL` alert is sent if the service in question is down or in a catastrophic state, such as a hard drive that is completely out of space and preventing the applications using that drive from running. An `OK` alert is sent when the service has recovered or become available again, such as when the total space used on the hard drive has dropped below the threshold set to alert for `CRITICAL` or `WARNING`.

OS and Hardware Checks

When monitoring any machine, it's best to start out with the checks at the OS and hardware layer. Even if you are running Cassandra in a virtualized environment such as Amazon or Rackspace, there are still hardware(ish) checks that should be instituted.

Disks and Partitions

The first thing you are going to want to check is the amount of free disk space on data partitions and the `CommitLog` partitions (assuming they are on separate partitions). Remember that if you are using `SizeTieredCompaction`, you shouldn't have the alert set for `WARNING` at 80% disk utilization and `CRITICAL` set at 90% disk utilization. The safer approach is to set the `WARNING` threshold to be roughly 35% disk utilization and the `CRITICAL` threshold at 45% disk utilization. `SizeTieredCompaction` is capable of taking up two times the size of the largest `SSTable` on disk. And while it is unlikely that a single `SSTable` would be 50% of the data on disk, it is better to be safe than sorry. Recovering from having too much data on disk is extremely difficult.

This concept of monitoring partitions and drives is also important because of `JBOD` support in Cassandra 1.2 and later. This means that Cassandra can have a single data directory on multiple disks. You will need to know if one or more of those disks are having an issue or require replacement. By monitoring the utilization and health of all the disks in your system, you will know their state and whether they need replacing or maintenance.

Last, you want to ensure that the drive that contains the log files doesn't fill up. Depending on your log settings, Cassandra has the potential to be very verbose in the log files. If the log files become too large, they can prevent the rest of your system from working if the drive(s) runs out of space.

Swap

Linux divides its physical memory into smaller chunks called pages. Swapping is the process whereby a page of memory is copied from memory to a dedicated space on the hard disk called swap space to free up that page of memory. Although there are cases where it is OK, it is normally not recommended for systems to be in a state where they are swapping memory. Typically, anything more than 5% to 10% of your swap space being used is cause for investigation.

On a Cassandra node, swapping is usually a bad sign, so you will want to monitor the swap partition for usage of nearly any kind. Since you should be able to hold the entire JVM's heap space in memory with at least a little room to spare for the operating system,

getting to the point of swapping out pages of memory means it might be a little too late to recover. One of the reasons Cassandra is able to function so well with regard to writes is the fact that many of the writes occur to the memory-mapped MemTables. Having these MemTables swap to disk would drastically impair the performance of Cassandra and should therefore be avoided when possible.

Clock Drift

Clock drift refers to the phenomenon where one clock does not run at the exact same speed as another clock. It is especially important to be aware of this if you are running in a virtualized environment as drift from the hypervisor can be much more prevalent than on regular iron. The system clock is incredibly important to Cassandra's write and reconciliation architecture. Most writes are serialized by timestamp. In other words, if two writes come in for the same column at almost the same time, the determining factor for which value wins is which timestamp is higher. If the system clocks in the ring are not all in sync, you are probably going to see some really strange behavior.

One of the ways to deal with that is to monitor the clock drift using NTP. NTP, or Network Time Protocol, is the most commonly used time synchronization system on the Internet. It also comes with a binary for telling you the offset (drift) from its synchronizing time server. You obviously want to minimize the amount of drift your system experiences. But there will invariably be some that you have to deal with. Monitoring is the way you know if the NTP daemon isn't doing the job it is supposed to be doing and keeping your clocks in sync. Being alerted to a problem with the clocks in a distributed environment that relies heavily on time for decision making could save a lot of time tracking down weird problems later on.

Ping Times

It is also a good idea to check the ping time responses from each of the Cassandra nodes being monitored. There are any number of reasons that these responses can begin to come back slowly. A few examples include the following:

- A machine that is doing too much work and running short of CPU cycles to respond quickly
- I/O saturation, too high an await (average wait) time, and the machine cannot respond quickly to the request
- Network saturation due to unthrottled streaming on a high-speed network link

Whatever the reason is, it is good to know if there is network congestion of which you should be aware. When a node is slow to receive packets (which is the case with nodes with high ping times), writes can be slow to come in and register, reads and writes will be dropped to keep up with the demand being put on the system, or any number of other weird behaviors may appear. What constitutes a high ping time from your monitoring server depends to a great extent on your network paths. Run a few ping tests from your monitoring server to your Cassandra nodes during regular usage periods to get a feel for what a normal threshold is.

CPU Usage

Cassandra is usually an I/O-bound system. You usually run into problems with disk writes or reads slowing down long before you run into CPU-related slowdown. But just to be safe, as different workloads call for different tools to be used at different times, you should monitor CPU usage. While there are many things you could look for when monitoring CPU usage, such as context switches or interrupt requests, a good place to start is usually watching the system load average. The system load average is an average of the number of processes waiting to get into the system's run queue over a period of time. In the case of the `uptime` command, it's over one, five, and 15 minutes. Keep in mind that in the case of multiprocessor systems, the load is relative to the number of processors and cores on the system.

The common rule for utilization is that you want to have a machine working hard but not overworking. This means that you typically want to have the machine running at about 70% utilization. That leaves you headroom for spikes in work and doesn't leave the machine underutilized during slower periods. So if you have four cores, having the load sit at around 3.00 is usually a safe bet. If you have four cores and the load is 3.5 or higher, you should try to find out what's wrong and fix it before things go from bad to worse.

Cassandra-Specific Health Checks

Once you have the basic system checks in place, it's time to add monitoring that is specific to Cassandra. There are various checks that interact with Cassandra at different levels of the system. Some are superficial such as checking to see if ports are alive and being listened on. Some checks require using a slightly more in-depth toolset to programmatically check the MBeans described earlier.

Ports

There are three primary ports of interest to Cassandra: 7000 (or 7001 if SSL/TLS is enabled), 7199, and 9160. Port 7000/7001 is used by Cassandra for cluster communication. This includes things such as the Gossip protocol and failure detection. Port 7199 is used by JMX. Port 9160 is the Thrift port and is used for client communication. In order for your cluster to function properly, all of these ports should be accessible.

While it is not necessary to specifically monitor these ports, it is a good idea to test them out one way or another. Testing the Thrift port (9160) is just testing whether you can connect to an instance using a Cassandra driver. In terms of monitoring, if you can connect, the check passes. If you can't connect to the server, the check should send off an alert. You can also use a simple TCP check here even though it is less comprehensive.

JMX Checks

Using some of the knowledge we gained from looking at the normal behavior of our system with JConsole, we are going to add some checks using JMX. There are plug-ins for Nagios that enable you to run JMX queries and compare the results against a set of predetermined thresholds. While there are many values that can be monitored through JMX, there are a few that stand out.

The first set of JMX checks to create is for read and write request latency. These values are given in microseconds because they should be that small. These latencies can be measured at the Cassandra application level and/or at the ColumnFamily level. Measuring them at the application level is important as a general health metric. High request latencies can be indicative of a bad disk or that your current read pattern is starting to slow down. If there is a ColumnFamily for which it is particularly important to have extremely low-latency reads and/or writes, it would be a good decision to monitor the performance for that ColumnFamily as well. It is important to note that read latency and write latency are two separate metrics provided by Cassandra, and both are important in their own right depending on your workload.

The next set of JMX metrics to keep tabs on is garbage collection timing. Cassandra will not only tell you how long its last garbage collection took but also how long that last ParNew GC took. A good way to think of ParNew garbage collection is that it is a stop-the-world garbage collection that uses multiple GC threads to complete its job. If you are monitoring the amount of time these take, you can easily set up an alert for when they start to take too long. Cassandra is unavailable during a stop-the-world garbage collection pause. The longer these pauses take, the longer Cassandra will be unavailable.

Another metric that is useful in helping to determine whether or not you need to add capacity to your cluster is PendingTasks under the CompactionManagerMBean. Depending on the speed and volume with which you ingest data, you will need to find a comfortable set of thresholds for your system. Typically, the number of PendingTasks should be relatively low, as in fewer than 50 at any given time. There are certainly acceptable reasons for things to back up, such as forced compactions or cleanup, but it is advisable to watch this metric carefully. If you have an alert set for PendingTasks and find this alert firing regularly, you may need to add more capacity (either more or faster disks or more nodes) to your cluster to keep up with the workload.

The last JMX metrics that should make it onto your first round of monitoring are the amount of on-heap and the amount of off-heap memory used at a time. The amount of on-heap memory used should always be less than the amount of heap that you have allowed the JVM to allocate. Since you know what this value is at start time, you should be able to easily monitor whether or not you are approaching that value. Off-heap memory tracking is a little harder to monitor for sane values. This is a metric where you will once again have to take a look at JConsole and see what regular and peak values are for the system under normal and peak operational loads so you don't send off useless alerts.

Log Monitoring

There is a lot of useful information in the Cassandra logs that can be indicative of a problem. As mentioned earlier in the chapter, you can find READ and WRITE dropped message counts within the INFO log level. There is a Nagios plug-in that can monitor logs and check for specific log messages. Using this plug-in, you can have Nagios alert you not just when there are READ and/or WRITE messages dropped, but you also can have it alert you when this happens more than *n* times per period. For instance, your application may be tolerant of missing READS and much less tolerant of missing WRITES. So the log monitoring check can alert you with a CRITICAL alert if more than 1,000 mutations

have been dropped over a five-minute period and with a `WARNING` alert if more than 1,000 mutations have been dropped over a 15-minute period.

This is just in the case of bad things happening in the `INFO` level. You can also have the log monitoring system alert you if any `FATAL`, `ERROR`, or `WARNING` log messages are put into the logs. Many of these plug-ins are configurable enough to send the log messages (or at least the one that caused the notification) along with the alert.

Cassandra Interactions

Now that we have the OS and system layer monitored and we know Cassandra is up and at least responding, it's time to check a little deeper. The further into the application you monitor, the better you will be able to sleep at night knowing things are functioning the way you want them to. Although it is useful and necessary to have superficial checks like load average and memory, the real value of monitoring systems is realized as you get deeper into the application.

What this means is that you should be checking things that are specific to your application in addition to the Cassandra server. If your application writes to a new `ColumnFamily` at the beginning of every month, you should have your monitoring system check before the month turnover that the new `ColumnFamily` exists (and optionally create it if it doesn't).

Another good use of monitoring resources is to check the response time of certain queries. If you are regularly running queries that roll up all the events for an hour, monitor how long that query takes to run and set up an alert if it's outside the normal threshold. In other words, if the query runs too fast, you want to know because it's possible you aren't collecting all the data you expect to be there. If the query takes too long to run, your system could be under heavy load or you may have just hit a point where you need to rethink your query patterns. Either way, that type of instrumentation is useful to measure how your system actually performs compared to how you expect it to perform.

If you run an application at the top of every hour—an extract, transform, load (ETL) process, for example—it might be a good idea to have the application put a “run complete” column somewhere when it's done. At the beginning of every hour, the monitoring system can run a query to check for the existence of the column for the last hour. If the “run complete” column doesn't exist for the last hour, it would be good to know so you can look into why.

Summary

There are many tools available for building monitoring systems. Nagios is just one of the common general-purpose monitoring tools. As long as some application is checking on the health and availability of your system and letting you know when an issue is present, or about to present itself, you will be in good shape. There are also some good examples of how your main application and other parts of your application interact with Cassandra and can be instrumented to give you a feeling of total information awareness and potentially the ability to get a good night's sleep when it is all in production.

In this example, Nagios can act as an early-warning tool. It can give you a heads-up to look at the machine in question and dig deeper into a potential problem before it turns into something more serious such as a full-fledged outage or a completely downed node. Ensuring an intelligently set-up monitoring infrastructure is essential to having a well-designed and architected system.

This page intentionally left blank

Index

A

- access.properties** file, for authentication/authorization, 128–129
- ACID (Atomicity, Consistency, Isolation, Durability)** database properties, 3
- active-active data centers**, 142
- Acunu Analytics**, 152–153
- ad hoc queries**, 140, 147
- ALL** option, for ColumnFamily tuning, 61
- ALLOW FILTERING** option, in CQL 3, 37
- ALTER KEYSPACE** command, in CQL 3, 31
- Amazon Web Services**, for running Cassandra, 50, 138
- analytics**
 - integrated, 154
 - low-latency, 152–153
 - real-time, 142
- Apache Cassandra**. *see* **Cassandra**
- Apache Hadoop**. *see* **Hadoop**
- approximate aggregates**, with Acunu, 152
- architecture**
 - peer-to-peer, 142
 - staged event-driven, 133–134
- archive_command** parameter, for CommitLog segments, 81
- asymmetrical replication**, 42–43
- atomicity property**, 3
- Atomicity, Consistency, Isolation, Durability (ACID)** database properties, 3
- authentication**, meta keyspaces and, 128–129

authorizer property, for Cassandra configuration, 13

availability

- with Cassandra, 142
- of transactions, 4

B

backups

- in Cassandra, 79
- using snapshots, 79–80

barriers=0 setting option, 58

BASE (Basically Available, Soft state, Eventual consistency) database properties, 4–5

Bash script, 82

basically available, as database property, 4–5

BATCH statement

- in CQL 3, 35–36
- updating counters using, 22

big-data techniques, with Pentaho, 154

BigTable data model, 135

binaries, installation from, 12

bloom filters

- for data structure accuracy, 61–62
- purpose/function of, 131–132
- SSTables and, 130–131

Brewer's theorem, 3–4

business intelligence solutions, with Pentaho, 154

ByteOrderedPartitioner, advantage of, 47

C

C# driver for Cassandra

- connecting to/disconnecting from cluster, 104–105
- creating sample class with, 104
- creating schema/writing data with, 105–106
- full C# sample class, 106–108

caching

- in Cassandra, 59
- ColumnFamily tuning, 61
- general tips for, 59–60
- global tuning, 60–61
- OOM errors and, 124

CAP (Consistency, Availability, Partition tolerance) theorem, 3–4

Cassandra

- applications, monitoring, 96
- C# driver for, 104–108
- caching in, 59
- current drivers for, 99
- data model, 17–19
- features of, 5–6
- for global data storage, 137–141
- health checks specific to, 94–96
- for high-volume real-time data, 141–147
- history of, 6
- Instaclustr managed hosting of, 154–155
- Java driver for, 100–104
- Python driver for, 108–112
- Ruby driver for, 112–116
- terminology, 8
- utilization of/choosing, 7
- for video analytics, 135–137

Cassandra Query Language. see CQL 1 (Cassandra Query Language 1); CQL 2 (Cassandra Query Language 2); CQL 3 (Cassandra Query Language 3)

cassandra.yaml file

- commitlog_directory in, 53
- for configuring Cassandra, 13
- snitches configured in, 43

CentOS, Cassandra installation from, 12

central processing unit. see CPU (central processing unit) usage

CL (consistency level)

- in ACID property, 3
- in CAP theorem, 4
- with Cassandra, 6
- choosing setting for, 147
- with reads/writes, 55–56
- specifying, 8–9
- tunable nature of, 147

cleanup, **with nodetool**, 75–76

ClientRequestMetrics, in MBeans, 91

clock drift, monitoring of, 93

cloud platforms, for running Cassandra, 49–50

cloud storage services, with Instaclustr, 154–155

cluster(s)

- balanced, 49
- as Cassandra term, 8
- connecting to, 100, 104–105, 112
- disconnecting from, 101, 105, 113
- Hadoop, 135
- multitenant/single-use, 138
- nodetool management of. *see* nodetool
- setup for single/multiple, 15–16

`cluster_name` **property**, in Cassandra configuration, 13

CLUSTERING KEYS, in data storage, 17–18

collections, in data modeling, 22–24

ColumnFamilies. *see also* specific

ColumnFamilies

- adjusting bloom filters for, 62
- caching within, 61, 124
- in Cassandra, 6
- compaction strategies for, 77
- compression settings for, 57
- counter, 22
- for customer records, 139
- information in System Keyspace, 127–128

- nodetool statistics on, 73–74
- schema-less, 7
- static/dynamic, 28
- taking snapshots of, 79–80
- wide-row, 17, 28, 136, 139–140

ColumnFamilyStatistics file, for data storage, 131

columns

- CQL 2 and, 28
- tombstones for deleted, 133

CommitLog directory(ies)

- archiving/restoring, 81–82
- Cassandra installation and, 11
- mutation operations and, 130
- optimizing, 53–54
- snapshots and, 79

`commitlog_directory` **property**, 14

`commitlog_segment_size_in_mb` **property**, 14

`commitlog_sync` **property**, 14

`commitlog_sync_period_in_ms` **property**, 14

compaction(s)

- large write workload and, 142, 147
- strategies for, 77
- types of, 76–77, 132
- unthrottling using nodetool, 77–78

CompactionManager, information in MBeans, 91

COMPOUND KEYS, data storage with, 17–19

compression

- benefits of, 141
- at ColumnFamily level, 57
- at network level, 56–57
- SnappyCompressor/DeflateCompressor for, 58

concurrency

- control of, 55
- SEDA model for, 133–134

`concurrent_reads` property, 15
`concurrent_writes` property, 15
ConcurrentLinkedHashMapProvider
 for efficient cache usage, 124
 global cache tuning and, 60–61
ConcurrentMarkSweep collector, 123
configuration, of Cassandra, 13–15
consistency. see CL (consistency level)
Consistency, Availability, Partition tolerance (CAP) theorem, 3–4
counter ColumnFamilies, creation of, 22
counter type, in CQL 3, 30
CPU (central processing unit) usage
 compression and, 141
 monitoring of, 94
CQL 1 (Cassandra Query Language 1), 27
CQL 2 (Cassandra Query Language 2), 28
CQL 3 (Cassandra Query Language 3)
 commands supported by, 30–37
 data types supported by, 28–29
 example schemas using, 37–39
 features of, 28
CREATE INDEX command, in CQL 3, 34
CREATE KEYSPACE command, in CQL 3, 31
CREATE TABLE/COLUMNFAMILY command, options for, 31–33
CRITICAL alert, with Nagios, 92
cross_node_timeout option, setting, 53
customer records, ColumnFamilies for, 139

D

dashboards, analytics, 152–153

data center(s)

 Cassandra clusters with multiple, 146
 number of replicas per, 42–43

data directories

 affecting CommitLog performance, 53
 Cassandra installation and, 11

data files, SSTables and, 131

data modeling, in Cassandra

 challenges with, 147
 collections in, 22–24
 overview of, 17–19
 query patterns for, 19–22
 sorting raw event data with, 144, 145

data partitions, Nagios monitoring, 92

data reading

 in C#, 106
 in Java, 102
 in Python, 110
 in Ruby, 114

data types, supported by CQL 3, 28–30

data visualization techniques, 152–153

data writing

 in C#, 105–106
 in Java, 101–102
 in Python, 110
 in Ruby, 113–114

data_property_directories property, in Cassandra configuration, 14

data=journal, commit=15 setting option, 58

data=writeback, nobh setting option, 58–59

DataStax Enterprise Cassandra, 151–152

date type, in CQL 3, 29–30

dateOf () function, in CQL 3, 30

DB (database)

 graph, 153
 NoSQL, 2, 142
 section of MBeans, 91
 transactions, ACID properties of, 3

DBMS (database management systems)

 Cassandra features as, 5–6
 distributed vs. relational, 2

Debian, installation from, 12

DEBUG logging level, 84

DeflateCompressor, 57, 58

DELETE command, in CQL 3, 35

deleting snapshots, 80–81

disk health, Nagios monitoring, 92

disk_failure_policy property, in Cassandra configuration, 14

drain, with nodetool, 75

DROP INDEX command, 34

DROP KEYSPACE command, 31

DROP TABLE command, 33

dstat tool, 120–121

durability

- performance tuning for, 55–56
- of transactions, 3

DynamicSnitch, keyspace creation and, 43–44

E

eBay, utilizing Cassandra, 141–147

EBS (Elastic Block Store) volumes, 50

EC2 (Elastic Computer Cloud), 50

Ec2MultiRegionSnitch, in keyspace creation, 44–45

Ec2Snitch, in keyspace creation, 44–45

edges, in distributed graph databases, 153

e-hail app, utilizing Cassandra, 137–141

EndPointSnitchInfo, in DB section of MBeans, 91

/etc/security/limits.conf settings, in performance tuning, 64

event attribute data, storing of, 136

event_metric value, creating tables for, 22

event_type value, storage of, 20

eventual consistency in database system, 5

ext4 file system, formatting devices for, 58–59

F

FailureDetector information, in MBeans, 91

false positives, with bloom filters, 61–62, 131–132

file limit settings, performance tuning and, 64

file system, for Cassandra deployment, 58–59

firewall ports, 49

flush, with nodetool, 75, 79

fraud detection, with Cassandra, 143–144

G

GC (garbage collection)

- JVM performance and, 66
- monitoring long-running, 123
- monitoring timing of, 95
- pauses in, 147
- viewing from Memory tab, 87, 88

GCGraceSeconds, for tombstones, 133

GCInspector, 123

global cache tuning, 60–61

Gossip protocol

- in cluster setup, 15
- detecting failure of nodes, 130
- information, in MBeans, 91
- purpose/utilization of, 129–130

Gossiper information, in MBeans, 91

graph database, Titan, 153

graph-based recommendation system, for taste profiles, 144–146

Gremlin graph traversal language, 153

H

Hadoop

- analytics with, 146, 151
- scalability of, 135, 136

Hailo taxi app, utilizing Cassandra, 137–141

health checks. *see* system health checks

high-volume real-time data, Cassandra capabilities for, 143–144

HintedHandoffManager, in DB section of MBeans, 91

HintedHandoffs, purpose/function of, 131

homogenous environment, as Cassandra term, 8

horizontal scalability, Cassandra for, 138

hot spots

- clustering order to remove, 21
- heavy read/write load leading to, 19

I

I/O (input/output)

- concurrency, testing of, 62–63
- increasing capacity, 122
- iostat monitoring, 119–120

idempotent operations, 140–141

index files, SSTables and, 130–131

IndexInfo ColumnFamily, 127

INFO logging level, 84

initial_token property

- in cluster setup, 15
- for configuring Cassandra, 13

INSERT command, in CQL 3, 34

Instaclustr, managed Cassandra hosting with, 154–155

installation process

- from Debian, RedHat/CentOS/Oracle, binaries, 11–12
- directories for, 11

insufficient resource errors, 124–126

integrated analytics platform, with Pentaho, 154

Internal section of MBeans, 91

internode communication

- information, in MBeans, 91
- Port 7000/7001 for, 49, 94

internode_compression controls, 56–57

iostat tool

- showing normal wait time for I/O, 119–120
- showing overly active system, 120

isolation property, affecting CommitLog performance, 3

J

Java driver for Cassandra

- connecting to/disconnecting from cluster with, 100–101
- creating sample class with, 100
- creating schema/writing data with, 101–102
- full Java sample class, 102–104
- reading data with, 102

JConsole

- logging in, 86–87
- MBeans tab in, 87, 90–91
- Memory/Threads tabs in, 87, 88–89

JMX (Java Management Extensions)

- features of, 85–86
- handshake, 49
- health checks, 94–95
- JConsole and, 86–91
- Port 7199 for, 49, 94

JVM (Java Virtual Machine)

- garbage collection and, 66, 123
- options for, 65
- for running Cassandra, 49
- setting maximum heap size for, 65–66

K

key cache

- global cache tuning and, 60
- size, OOM errors and, 124

key/value stores, 5–6, 7

KEYS_ONLY option, for ColumnFamily tuning, 61

keyspace creation

- firewalls in, 49
- node layout in, 48–49
- overview of, 41
- partitioners and, 46–48
- platforms in, 49–50
- replication strategies in, 41–43
- snitches and, 43–46

keyspaces, meta, 127–129

L

LeveledCompactionStrategy, 77, 132, 137

LIMIT option, in CQL 3, 37

limits.conf file, 125

linear scalability, with Cassandra, 142

listen_address property, in Cassandra configuration, 15

lists, for data model collections, 23–24

log monitoring, 95–96

logging levels

- changing, 84
- mutation/dropped READ messages in, 84–85
- overview of, 83

low-volume application, data model for log storage in, 20–21

M

major compactions, 76–77

maps, for data model collections, 24

MAX_HEAP_SIZE value, JVM performance and, 65–66

maxTimeuuid function, in CQL 3, 30

MBeans (Managed Beans) tab, with JConsole, 86, 87, 90–91

memory

- heap usage of, 87, 88

- monitoring swapping of, 92–93
- on-heap/off-heap usage, 95, 147
- swap setting and, 63

Memory tab, with JConsole, 87, 88

memtable_total_space_in_mb property, for configuring Cassandra, 15

MemTables

- caching affect on, 60
- flushing data from memory, 75
- mutation operations and, 130
- OOM errors and, 124
- performance tuning of, 54–55

MessagingService information, in MBeans, 91

meta keyspaces

- for authentication, 128–129
- overview of, 127
- System Keyspace, 127–128

Metrics section of MBeans, 91

Migrations ColumnFamily, 127–128

minor compactions, 76–77

minTimeuuid function, in CQL 3, 30

mobile notification tracking, with Cassandra, 143

multitenant clusters, 138, 146

Murmur3Partitioners, 48

mutation operations, CommitLogs/MemTables and, 130

N

Nagios

- clock drift/ping times and, 93
- CPU usage and, 94
- monitoring disks/partitions/drives, 92
- monitoring swap partition, 92–93
- primary alerts in, 91–92

naming conventions, for snapshots, 80

Net section of MBeans, 91

network compression, 56–57

Network Time Protocol (NTP), monitoring
clock drift, 93

NetworkTopologyStrategy, for replication,
42–43

`noatime` **setting option**, 58

node(s)
as Cassandra term, 8
communication between, 49
detecting failure of, 130
freezing, troubleshooting for, 123
information about, 72
ring view differing between, 124
seed, 15
virtual, 48–49

`nodetool`
cleaning with, 75–76
in cluster setup, 16
ColumnFamily statistics with, 73–74
common commands, 121
flushing/draining with, 75
function of, 69
general usage of, 71–72
node information with, 72
options for, 69–71
taking snapshots with, 80–81
thread pool statistics with, 74–75
three-node cluster information,
72–73
unthrottling compactions with, 77–78

`nodetool cfstats`, 122

`nodetool info` **output**, 72

`nodetool repair`, 76

`nodetool ring` **command**, 72–73

`nodetool scrub`, 76

`nodetool upgradesstables`, 76

`nodetool version` **output**, 72

NONE **option**, for ColumnFamily tuning, 61

NoSQL databases

active-active data centers in, 142
overview of, 2

`now()` **function**, in CQL 3, 30

NTP (Network Time Protocol), monitoring
clock drift, 93

`num_tokens` **property**, for configuring
Cassandra, 13

O

OOM (out-of-memory) errors, tracking of, 124

Ooyala online video analytics, utilizing
Cassandra, 135–137

OpsCenter, DataStax, 141, 151–152

Oracle, Cassandra installation from, 12

order and shipment tracking, with Cassandra,
143–144

ORDER BY **option**, in CQL 3, 36

P

ParNew collector, 123

partition tolerance, of transactions, 4

`partitioner` **property**, in Cassandra
configuration, 14

partitioners
ByteOrderedPartitioners, 47
function of, 46–47
Random/Murmur3Partitioners, 47–48

`password.properties` **file**, for authentication/
authorization, 128–129

peer-to-peer architecture, 142

PendingTasks, monitoring performance of, 95

Pentaho integrated analytics platform, 154

performance tuning
adjusting MemTables, 54–55
bloom filters for, 61–62
caching in, 59–61
compression and, 56–58

concurrency in, 55
 for durability/consistency, 55–56
 file system for, 58–59
 JVM tuning for, 65–66
 memory/swap setting in, 63
 methodology for, 51–52
 optimizing CommitLog, 53–54
 setting timeouts, 52–53
 solid-state drives in, 64–65
`sysctl` network/file limit settings in, 63
 testing I/O concurrency, 62–63
 testing in production, 52
permissions_validity_in_ms property,
 in **Cassandra** configuration, 14
ping time responses, monitoring of, 93
platforms, for running Cassandra, 49–50
plug-and-play capabilities, of Acunu, 152
port(s)
 default JMX, 86–87
 for internode communication, 49
 monitoring health of, 94
PRIMARY KEY operator
 CQL 3 and, 28
 CREATE TABLE command and, 31–32
 data storage with, 17–18
**PropertyFileSnitch, keyspace creation and,
 45–46**
Python driver for Cassandra
 connecting to/disconnecting from
 cluster/creating schema with, 109
 creating sample class with, 108
 full Python sample class, 110–112
 writing data/reading data with, 110

Q

queries
 pre-aggregating/grouping of, 152
 prior identification of, 140, 147

query patterns
 counter ColumnFamilies and, 22
 for low-volume application, 20–21
 optimized, 21
 for relational database, 19–20
**quorum read/writes, defining consistency
 and, 8**

R

RackInferingSnitch, in keyspace creation, 44
RAID0, for running Cassandra, 49–50
RandomPartitioners, types of, 47–48
Raw Event Data ColumnFamily, 144, 145
 raw event data, storage of, 136
**RDBMSs (relational database management
 systems)**
 data model for log storage in, 19–20
 as differing from **Cassandra**, 6, 7
read latency
 LeveledCompaction and, 132
 monitoring of, 95
 strict requirements for, 142
 troubleshooting, 122
**real-time analytics, with Cassandra, 142,
 143–144**
RedHat, Cassandra installation from, 12
relational database management systems.
 see **RDBMSs (relational database
 management systems)**
remote procedure call (RPC) framework, 27
replica(s)
 counts, 42–43
 partitioner placement of, 46–48
replication factor. see RF (replication factor)
replication strategies
 multi-data-center, 6
 NetworkTopologyStrategy, 42–43
 SimpleStrategy, 42

Request section of MBeans, 91

`restore_command`, for **CommitLog**
 archiving, 81

`restore_directories` **parameter**, for
 archived CommitLogs, 81

`restore_point_in_time` **parameter**, for
 archived CommitLogs, 82

RF (replication factor), 8
 choice of setting for, 147
 definition of/setting, 41
 for resilient data storage, 138

ring view, differing between nodes, 124

row cache, **global cache tuning** and,
 60–61

row stores, **Cassandra features** of, 6

rows
 caching and, 59
 CQL 2 and, 28

ROWS_ONLY option, for **ColumnFamily tuning**,
 61

RPC (remote procedure call) framework, 27

Ruby driver for Cassandra
 connecting to/disconnecting from
 cluster with, 112–113
 creating sample class with, 112
 creating schema with, 113
 full Ruby sample class, 115–116
 writing data/reading data with,
 113–114

S

`saved_caches_directory` **property**, in
Cassandra configuration, 14

scalability factor, with **Cassandra**, 138, 142

schema
 creation of, 101, 105, 109, 113
 migrations of, 127–128
 option for creating, 7
 time-series wide row, 136

SEDA (staged event-driven architecture), for
concurrency, 133–134

`seed_provider` **property**, in **Cassandra**
configuration, 14–15

SELECT statement, in **CQL 3**, 36

SerializingCacheProvider, 60–61, 124

sets, for data model collections, 22–23

sharding, built-in, 142

SimpleAuthenticator **setting**, 128

SimpleSnitch, **keyspace creation** and, 43

SimpleStrategy, for replication, 42

single-use clusters, 138

SizeTieredCompactionStrategy, 77, 132

SnappyCompressor, 57, 58

snapshots
 function of, 79
 removing, 80–81
 taking/naming, 79–80

snitches
 definition of/**SimpleSnitch**, 43
 DynamicSnitch, 43–44
 Ec2MultiRegionSnitch, 44–45
 Ec2Snitch, 44–45
 PropertyFileSnitch, 45–46
 RackInferringSnitch, 44

soft state database system, 5

Solr search platform, 151

Spark computing framework, 136

SQL (Structured Query Language), **Cassandra**
 and, 27, 140

SSDs (solid-state drives), **tuning** of,
 64–65

SSH port (22), 49

SSTable(s)
 count, monitoring, 122
 function of files in, 130
 nodetool rebuilding, 76

sstables scrub tool, 76

staged event-driven architecture (SEDA), for concurrency, 133–134

streaming_socket_timeout_in_ms value, 53

StreamingService information, in MBeans, 91

Structured Query Language (SQL), Cassandra and, 27, 140

swap setting, virtual memory and, 63

swapping memory

- long-running GCs and, 123
- monitoring, 92–93

sysctl network settings

- performance tuning and, 64
- updating for max_map_count, 126

system health checks

- Cassandra-specific, 94–96
- with Nagios, 91–94

System Keyspace, 127–128

T

tables

- changing cache setting on, 61
- creating static/dynamic, 32
- options for creating, 32–33
- PRIMARY KEY and, 31–32

Taste Graph modeled in Cassandra, 144–146

taxi app, utilizing Cassandra, 137–141

terminology, Cassandra, 8

text search capability, with Cassandra, 151

thread pools

- for SEDA stages, 133–134
- statistics on with nodetool, 74–75

thread stack size, JVM performance and, 65

Threads tab, with JConsole, 87, 89

three-node ring

- compaction and, 77
- nodetool information on, 72–73

Thrift port (9160), 49, 94

Thrift RPC framework, 27

timeouts, configuration of, 52–53

time-series data

- Cassandra handling, 5
- eBay and, 143–144
- Hailo and, 139–140
- Ooyala and, 136

TimeUUID types, in CQL 3, 30

Titan distributed graph database, 153

token(s)

- in cluster setup, 15–16
- nodetool cleanup and, 76
- ranges, vnodes and, 48–49

tombstones, function of, 132–133

tools, troubleshooting

- dstat/nodetool, 120–121
- iostat, 119–120

traversing, in distributed graph databases, 153

troubleshooting

- insufficient resource errors, 124–126
- long-running GCs, 123
- OOM errors/differing ring view, 124
- slow reads/fast writes, 122
- tools for, 119–121

TRUNCATE command, in CQL 3, 33–34

tunable consistency

- advantage of, 147
- as Cassandra term, 8–9

U

ulimit -a command, 125

unixTimestampOf() function, in CQL 3, 30

UPDATE command, 35

USE command, 31

V

Versions ColumnFamily, 128

vertices, in distributed graph databases, 153

video analytics aggregates, Cassandra handling, 135–137

virtual memory, swap setting and, 63

vnodes (virtual nodes), data distribution and, 48–49

W

WARNING alert, with Nagios, 92

Web scale technologies, 2

WHERE clause, in CQL 3, 36

wide rows

- CQL 2 for, 28

- data storage with, 17
- for time-series data, 136, 139–140

write latency, monitoring of, 95

X

XFS file system, 137

- XX:+CMSParallelRemarkEnabled **setting**, 66
- XX:+UseCMSInitiatingOccupancyOnly **setting**, 66
- XX:+UseConcMarkSweepGC **setting**, 66
- XX:+UserParNewGC **setting**, 66
- XX:CMSInitiatingOccupancyFraction=75 **setting**, 66