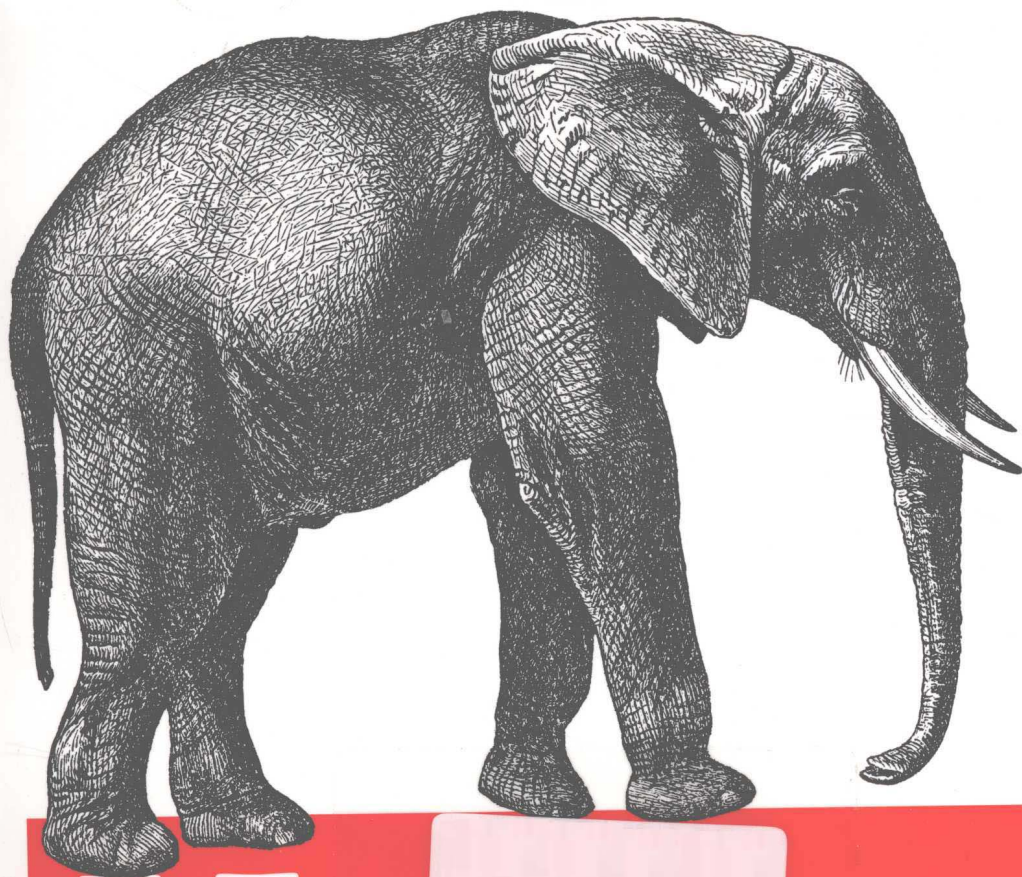


Hadoop权威指南 (影印版)

第三版
修订版



Hadoop

The Definitive Guide

O'REILLY®

东南大学出版社

Tom White 著

第三版

Hadoop权威指南 (影印版)

Hadoop: The Definitive Guide

Tom White

O'REILLY®

Beijing · Cambridge · Farnham · Köln · Sebastopol · Tokyo

O'Reilly Media, Inc. 授权东南大学出版社出版

东南大学出版社

图书在版编目 (CIP) 数据

Hadoop 权威指南: 第3版: 英文/(美)怀特 (White, T.)
著. —影印本. —南京: 东南大学出版社, 2013.1

书名原文: Hadoop: The Definitive Guide, 3E

ISBN 978-7-5641-3893-6

I. ① H… II. ① 怀… III. ① 数据处理—应用软件—
英文 IV. ① TP274

中国版本图书馆 CIP 数据核字 (2012) 第 273580 号

江苏省版权局著作权合同登记

图字: 10-2012-168 号

©2012 by O'Reilly Media, Inc.

Reprint of the English Edition, jointly published by O'Reilly Media, Inc. and Southeast University Press, 2013. Authorized reprint of the original English edition, 2012 O'Reilly Media, Inc., the owner of all rights to publish and sell the same.

All rights reserved including the rights of reproduction in whole or in part in any form.

英文原版由 O'Reilly Media, Inc. 出版 2012。

英文影印版由东南大学出版社出版 2013。此影印版的出版和销售得到出版权和销售权的所有者——O'Reilly Media, Inc. 的许可。

版权所有, 未得书面许可, 本书的任何部分和全部不得以任何形式重制。

Hadoop 权威指南 第三版 (影印版)

出版发行: 东南大学出版社

地 址: 南京四牌楼 2 号 邮编: 210096

出 版 人: 江建中

网 址: <http://www.seupress.com>

电子邮件: press@seupress.com

印 刷: 扬中市印刷有限公司

开 本: 787 毫米 × 980 毫米 16 开本

印 张: 43

字 数: 842 千字

版 次: 2013 年 1 月第 1 版

印 次: 2013 年 1 月第 1 次印刷

书 号: ISBN 978-7-5641-3893-6

定 价: 98.00 元 (册)

本社图书若有印装质量问题, 请直接与营销部联系。电话 (传真): 025-83791830

For Eliane, Emilia, and Lottie

Foreword

Hadoop got its start in Nutch. A few of us were attempting to build an open source web search engine and having trouble managing computations running on even a handful of computers. Once Google published its GFS and MapReduce papers, the route became clear. They'd devised systems to solve precisely the problems we were having with Nutch. So we started, two of us, half-time, to try to re-create these systems as a part of Nutch.

We managed to get Nutch limping along on 20 machines, but it soon became clear that to handle the Web's massive scale, we'd need to run it on thousands of machines and, moreover, that the job was bigger than two half-time developers could handle.

Around that time, Yahoo! got interested, and quickly put together a team that I joined. We split off the distributed computing part of Nutch, naming it Hadoop. With the help of Yahoo!, Hadoop soon grew into a technology that could truly scale to the Web.

In 2006, Tom White started contributing to Hadoop. I already knew Tom through an excellent article he'd written about Nutch, so I knew he could present complex ideas in clear prose. I soon learned that he could also develop software that was as pleasant to read as his prose.

From the beginning, Tom's contributions to Hadoop showed his concern for users and for the project. Unlike most open source contributors, Tom is not primarily interested in tweaking the system to better meet his own needs, but rather in making it easier for anyone to use.

Initially, Tom specialized in making Hadoop run well on Amazon's EC2 and S3 services. Then he moved on to tackle a wide variety of problems, including improving the MapReduce APIs, enhancing the website, and devising an object serialization framework. In all cases, Tom presented his ideas precisely. In short order, Tom earned the role of Hadoop committer and soon thereafter became a member of the Hadoop Project Management Committee.

Tom is now a respected senior member of the Hadoop developer community. Though he's an expert in many technical corners of the project, his specialty is making Hadoop easier to use and understand.

Given this, I was very pleased when I learned that Tom intended to write a book about Hadoop. Who could be better qualified? Now you have the opportunity to learn about Hadoop from a master—not only of the technology, but also of common sense and plain talk.

—Doug Cutting
Shed in the Yard, California

Preface

Martin Gardner, the mathematics and science writer, once said in an interview:

Beyond calculus, I am lost. That was the secret of my column's success. It took me so long to understand what I was writing about that I knew how to write in a way most readers would understand.¹

In many ways, this is how I feel about Hadoop. Its inner workings are complex, resting as they do on a mixture of distributed systems theory, practical engineering, and common sense. And to the uninitiated, Hadoop can appear alien.

But it doesn't need to be like this. Stripped to its core, the tools that Hadoop provides for building distributed systems—for data storage, data analysis, and coordination—are simple. If there's a common theme, it is about raising the level of abstraction—to create building blocks for programmers who just happen to have lots of data to store, or lots of data to analyze, or lots of machines to coordinate, and who don't have the time, the skill, or the inclination to become distributed systems experts to build the infrastructure to handle it.

With such a simple and generally applicable feature set, it seemed obvious to me when I started using it that Hadoop deserved to be widely used. However, at the time (in early 2006), setting up, configuring, and writing programs to use Hadoop was an art. Things have certainly improved since then: there is more documentation, there are more examples, and there are thriving mailing lists to go to when you have questions. And yet the biggest hurdle for newcomers is understanding what this technology is capable of, where it excels, and how to use it. That is why I wrote this book.

The Apache Hadoop community has come a long way. Over the course of three years, the Hadoop project has blossomed and spun off half a dozen subprojects. In this time, the software has made great leaps in performance, reliability, scalability, and manageability. To gain even wider adoption, however, I believe we need to make Hadoop even easier to use. This will involve writing more tools; integrating with more systems; and

1. "The science of fun," Alex Bellos, *The Guardian*, May 31, 2008, <http://www.guardian.co.uk/science/2008/may/31/maths.science>.

writing new, improved APIs. I'm looking forward to being a part of this, and I hope this book will encourage and enable others to do so, too.

Administrative Notes

During discussion of a particular Java class in the text, I often omit its package name to reduce clutter. If you need to know which package a class is in, you can easily look it up in Hadoop's Java API documentation for the relevant subproject, linked to from the Apache Hadoop home page at <http://hadoop.apache.org/>. Or if you're using an IDE, it can help using its auto-complete mechanism.

Similarly, although it deviates from usual style guidelines, program listings that import multiple classes from the same package may use the asterisk wildcard character to save space (for example, `import org.apache.hadoop.io.*`).

The sample programs in this book are available for download from the website that accompanies this book: <http://www.hadoopbook.com/>. You will also find instructions there for obtaining the datasets that are used in examples throughout the book, as well as further notes for running the programs in the book, and links to updates, additional resources, and my blog.

What's in This Book?

The rest of this book is organized as follows. Chapter 1 emphasizes the need for Hadoop and sketches the history of the project. Chapter 2 provides an introduction to MapReduce. Chapter 3 looks at Hadoop filesystems, and in particular HDFS, in depth. Chapter 4 covers the fundamentals of I/O in Hadoop: data integrity, compression, serialization, and file-based data structures.

The next four chapters cover MapReduce in depth. Chapter 5 goes through the practical steps needed to develop a MapReduce application. Chapter 6 looks at how MapReduce is implemented in Hadoop, from the point of view of a user. Chapter 7 is about the MapReduce programming model and the various data formats that MapReduce can work with. Chapter 8 is on advanced MapReduce topics, including sorting and joining data.

Chapters 9 and 10 are for Hadoop administrators and describe how to set up and maintain a Hadoop cluster running HDFS and MapReduce.

Later chapters are dedicated to projects that build on Hadoop or are related to it. Chapters 11 and 12 present Pig and Hive, which are analytics platforms built on HDFS and MapReduce, whereas Chapters 13, 14, and 15 cover HBase, ZooKeeper, and Sqoop, respectively.

Finally, Chapter 16 is a collection of case studies contributed by members of the Apache Hadoop community.

What's New in the Second Edition?

The second edition has two new chapters on Hive and Sqoop (Chapters 12 and 15), a new section covering Avro (in Chapter 4), an introduction to the new security features in Hadoop (in Chapter 9), and a new case study on analyzing massive network graphs using Hadoop (in Chapter 16).

This edition continues to describe the 0.20 release series of Apache Hadoop because this was the latest stable release at the time of writing. New features from later releases are occasionally mentioned in the text, however, with reference to the version that they were introduced in.

What's New in the Third Edition?

The third edition covers the 1.x (formerly 0.20) release series of Apache Hadoop, as well as the newer 0.22 and 2.x (formerly 0.23) series. With a few exceptions, which are noted in the text, all the examples in this book run against these versions. The features in each release series are described at a high level in “Hadoop Releases” on page 13.

This edition uses the new MapReduce API for most of the examples. Because the old API is still in widespread use, it continues to be discussed in the text alongside the new API, and the equivalent code using the old API can be found on the book's website.

The major change in Hadoop 2.0 is the new MapReduce runtime, MapReduce 2, which is built on a new distributed resource management system called YARN. This edition includes new sections covering MapReduce on YARN: how it works (Chapter 6) and how to run it (Chapter 9).

There is more MapReduce material, too, including development practices such as packaging MapReduce jobs with Maven, setting the user's Java classpath, and writing tests with MRUnit (all in Chapter 5); and more depth on features such as output committers, the distributed cache (both in Chapter 8), and task memory monitoring (Chapter 9). There is a new section on writing MapReduce jobs to process Avro data (Chapter 4), and one on running a simple MapReduce workflow in Oozie (Chapter 5).

The chapter on HDFS (Chapter 3) now has introductions to high availability, federation, and the new WebHDFS and HttpFS filesystems.

The chapters on Pig, Hive, Sqoop, and ZooKeeper have all been expanded to cover the new features and changes in their latest releases.

In addition, numerous corrections and improvements have been made throughout the book.

Conventions Used in This Book

The following typographical conventions are used in this book:

Italic

Indicates new terms, URLs, email addresses, filenames, and file extensions.

Constant width

Used for program listings, as well as within paragraphs to refer to program elements such as variable or function names, databases, data types, environment variables, statements, and keywords.

Constant width bold

Shows commands or other text that should be typed literally by the user.

Constant width italic

Shows text that should be replaced with user-supplied values or by values determined by context.



This icon signifies a tip, suggestion, or general note.



This icon indicates a warning or caution.

Using Code Examples

This book is here to help you get your job done. In general, you may use the code in this book in your programs and documentation. You do not need to contact us for permission unless you're reproducing a significant portion of the code. For example, writing a program that uses several chunks of code from this book does not require permission. Selling or distributing a CD-ROM of examples from O'Reilly books does require permission. Answering a question by citing this book and quoting example code does not require permission. Incorporating a significant amount of example code from this book into your product's documentation does require permission.

We appreciate, but do not require, attribution. An attribution usually includes the title, author, publisher, and ISBN. For example: "*Hadoop: The Definitive Guide*, Third Edition, by Tom White. Copyright 2011 Tom White, 978-1-449-31152-0."

If you feel your use of code examples falls outside fair use or the permission given here, feel free to contact us at permissions@oreilly.com.

Safari® Books Online



Safari Books Online (www.safaribooksonline.com) is an on-demand digital library that delivers expert content in both book and video form from the world's leading authors in technology and business. Technology professionals, software developers, web designers, and business and creative professionals use Safari Books Online as their primary resource for research, problem solving, learning, and certification training.

Safari Books Online offers a range of product mixes and pricing programs for organizations, government agencies, and individuals. Subscribers have access to thousands of books, training videos, and prepublication manuscripts in one fully searchable database from publishers like O'Reilly Media, Prentice Hall Professional, Addison-Wesley Professional, Microsoft Press, Sams, Que, Peachpit Press, Focal Press, Cisco Press, John Wiley & Sons, Syngress, Morgan Kaufmann, IBM Redbooks, Packt, Adobe Press, FT Press, Apress, Manning, New Riders, McGraw-Hill, Jones & Bartlett, Course Technology, and dozens more. For more information about Safari Books Online, please visit us online.

How to Contact Us

Please address comments and questions concerning this book to the publisher:

O'Reilly Media, Inc.
1005 Gravenstein Highway North
Sebastopol, CA 95472
800-998-9938 (in the United States or Canada)
707-829-0515 (international or local)
707-829-0104 (fax)

We have a web page for this book, where we list errata, examples, and any additional information. You can access this page at:

<http://www.oreilly.com/catalog/9781449311520>

To comment or ask technical questions about this book, send email to:

bookquestions@oreilly.com

For more information about our books, courses, conferences, and news, see our website at <http://www.oreilly.com>.

Find us on Facebook: <http://facebook.com/oreilly>

Follow us on Twitter: <http://twitter.com/oreillymedia>

Watch us on YouTube: <http://www.youtube.com/oreillymedia>

Acknowledgments

I have relied on many people, both directly and indirectly, in writing this book. I would like to thank the Hadoop community, from whom I have learned, and continue to learn, a great deal.

In particular, I would like to thank Michael Stack and Jonathan Gray for writing the chapter on HBase. Thanks also go to Adrian Woodhead, Marc de Palol, Joydeep Sen Sarma, Ashish Thusoo, Andrzej Bialecki, Stu Hood, Chris K. Wensel, and Owen O'Malley for contributing case studies for Chapter 16.

I would like to thank the following reviewers who contributed many helpful suggestions and improvements to my drafts: Raghu Angadi, Matt Biddulph, Christophe Bisciglia, Ryan Cox, Devaraj Das, Alex Dorman, Chris Douglas, Alan Gates, Lars George, Patrick Hunt, Aaron Kimball, Peter Krey, Hairong Kuang, Simon Maxen, Olga Natkovich, Benjamin Reed, Konstantin Shvachko, Allen Wittenauer, Matei Zaharia, and Philip Zeyliger. Ajay Anand kept the review process flowing smoothly. Philip (“flip”) Kromer kindly helped me with the NCDC weather dataset featured in the examples in this book. Special thanks to Owen O'Malley and Arun C. Murthy for explaining the intricacies of the MapReduce shuffle to me. Any errors that remain are, of course, to be laid at my door.

For the second edition, I owe a debt of gratitude for the detailed review and feedback from Jeff Bean, Doug Cutting, Glynn Durham, Alan Gates, Jeff Hammerbacher, Alex Kozlov, Ken Krugler, Jimmy Lin, Todd Lipcon, Sarah Sproehle, Vinithra Varadhara-jan, and Ian Wrigley, as well as all the readers who submitted errata for the first edition. I would also like to thank Aaron Kimball for contributing the chapter on Sqoop, and Philip (“flip”) Kromer for the case study on graph processing.

For the third edition, thanks go to Alejandro Abdelnur, Eva Andreasson, Eli Collins, Doug Cutting, Patrick Hunt, Aaron Kimball, Aaron T. Myers, Brock Noland, Arvind Prabhakar, Ahmed Radwan, and Tom Wheeler for their feedback and suggestions. Rob Weltman kindly gave very detailed feedback for the whole book, which greatly improved the final manuscript. Thanks also go to all the readers who submitted errata for the second edition.

I am particularly grateful to Doug Cutting for his encouragement, support, and friendship, and for contributing the Foreword.

Thanks also go to the many others with whom I have had conversations or email discussions over the course of writing the book.

Halfway through writing this book, I joined Cloudera, and I want to thank my colleagues for being incredibly supportive in allowing me the time to write and to get it finished promptly.

I am grateful to my editor, Mike Loukides, and his colleagues at O'Reilly for their help in the preparation of this book. Mike has been there throughout to answer my questions, to read my first drafts, and to keep me on schedule.

Finally, the writing of this book has been a great deal of work, and I couldn't have done it without the constant support of my family. My wife, Eliane, not only kept the home going, but also stepped in to help review, edit, and chase case studies. My daughters, Emilia and Lottie, have been very understanding, and I'm looking forward to spending lots more time with all of them.

Table of Contents

Foreword	xv
Preface	xvii
1. Meet Hadoop	1
Data!	1
Data Storage and Analysis	3
Comparison with Other Systems	4
Rational Database Management System	4
Grid Computing	6
Volunteer Computing	8
A Brief History of Hadoop	9
Apache Hadoop and the Hadoop Ecosystem	12
Hadoop Releases	13
What's Covered in This Book	15
Compatibility	15
2. MapReduce	17
A Weather Dataset	17
Data Format	17
Analyzing the Data with Unix Tools	19
Analyzing the Data with Hadoop	20
Map and Reduce	20
Java MapReduce	22
Scaling Out	30
Data Flow	30
Combiner Functions	33
Running a Distributed MapReduce Job	36
Hadoop Streaming	36
Ruby	36
Python	39

Hadoop Pipes	40
Compiling and Running	41
3. The Hadoop Distributed Filesystem	43
The Design of HDFS	43
HDFS Concepts	45
Blocks	45
Namenodes and Datanodes	46
HDFS Federation	47
HDFS High-Availability	48
The Command-Line Interface	49
Basic Filesystem Operations	50
Hadoop Filesystems	52
Interfaces	53
The Java Interface	55
Reading Data from a Hadoop URL	55
Reading Data Using the FileSystem API	57
Writing Data	60
Directories	62
Querying the Filesystem	62
Deleting Data	67
Data Flow	67
Anatomy of a File Read	67
Anatomy of a File Write	70
Coherency Model	72
Data Ingest with Flume and Sqoop	74
Parallel Copying with distcp	75
Keeping an HDFS Cluster Balanced	76
Hadoop Archives	77
Using Hadoop Archives	77
Limitations	79
4. Hadoop I/O	81
Data Integrity	81
Data Integrity in HDFS	81
LocalFileSystem	82
ChecksumFileSystem	83
Compression	83
Codecs	85
Compression and Input Splits	89
Using Compression in MapReduce	90
Serialization	93
The Writable Interface	94

Writable Classes	96
Implementing a Custom Writable	103
Serialization Frameworks	108
Avro	110
Avro Data Types and Schemas	111
In-Memory Serialization and Deserialization	114
Avro Datafiles	117
Interoperability	118
Schema Resolution	121
Sort Order	123
Avro MapReduce	124
Sorting Using Avro MapReduce	128
Avro MapReduce in Other Languages	130
File-Based Data Structures	130
SequenceFile	130
MapFile	137
5. Developing a MapReduce Application	143
The Configuration API	144
Combining Resources	145
Variable Expansion	146
Setting Up the Development Environment	146
Managing Configuration	148
GenericOptionsParser, Tool, and ToolRunner	150
Writing a Unit Test with MRUnit	154
Mapper	154
Reducer	156
Running Locally on Test Data	157
Running a Job in a Local Job Runner	157
Testing the Driver	160
Running on a Cluster	161
Packaging a Job	162
Launching a Job	163
The MapReduce Web UI	165
Retrieving the Results	168
Debugging a Job	170
Hadoop Logs	175
Remote Debugging	177
Tuning a Job	178
Profiling Tasks	179
MapReduce Workflows	181
Decomposing a Problem into MapReduce Jobs	181
JobControl	183

Apache Oozie	183
6. How MapReduce Works	189
Anatomy of a MapReduce Job Run	189
Classic MapReduce (MapReduce 1)	190
YARN (MapReduce 2)	196
Failures	202
Failures in Classic MapReduce	202
Failures in YARN	204
Job Scheduling	206
The Fair Scheduler	207
The Capacity Scheduler	207
Shuffle and Sort	208
The Map Side	208
The Reduce Side	210
Configuration Tuning	211
Task Execution	214
The Task Execution Environment	215
Speculative Execution	215
Output Committers	217
Task JVM Reuse	219
Skipping Bad Records	220
7. MapReduce Types and Formats	223
MapReduce Types	223
The Default MapReduce Job	227
Input Formats	234
Input Splits and Records	234
Text Input	245
Binary Input	249
Multiple Inputs	250
Database Input (and Output)	251
Output Formats	251
Text Output	252
Binary Output	253
Multiple Outputs	253
Lazy Output	257
Database Output	258
8. MapReduce Features	259
Counters	259
Built-in Counters	259
User-Defined Java Counters	264