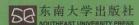


# Hadoop MapReduce v2 参考手册 第2版 (影印版)

adoop MapReduce v2 Cookbook econd Edition

Thilina Gunarathne 著

PACKT



# Hadoop MapReduce v2 参考手册

第2版 (影印版)

Thilina Gunarathne 著

南京 东南大学出版社

#### 图书在版编目(CIP)数据

Hadoop MapReduce v2 参考手册:第2版:英文/ (美)冈纳拉森(Gunarathne,T.)著.一影印本. 一南京:东南大学出版社,2016.1

书名原文: Hapdoop MapReduce v2 Cookbook, Second Edition

ISBN 978 - 7 - 5641 - 6089 - 0

Ⅱ.①H… Ⅲ.①冈… Ⅲ.①数据处理软件—手册—英文 ②软件工具—程序设计—手册—英文 Ⅳ.① TP274-62 ②TP311.56-62

中国版本图书馆 CIP 数据核字(2015)第 256601 号

#### © 2015 by PACKT Publishing Ltd

Reprint of the English Edition, jointly published by PACKT Publishing Ltd and Southeast University Press, 2016. Authorized reprint of the original English edition, 2015 PACKT Publishing Ltd, 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 PACKT Publishing Ltd 出版 2015。

英文影印版由东南大学出版社出版 2016。此影印版的出版和销售得到出版权和销售权的所有者—— PACKT Publishing Ltd 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 MapReduce v2 参考手册 第 2 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:常州市武进第三印刷有限公司

开 本: 787 毫米×980 毫米 16 开本

印 张: 20

字 数: 392 千字

版 次: 2016年1月第1版

印 次: 2016年1月第1次印刷

书 号: ISBN 978-7-5641-6089-0

定 价: 64.00元

## **Credits**

Authors

Thilina Gunarathne

Srinath Perera

Reviewers

Skanda Bhargav

Randal Scott King

**Dmitry Spikhalskiy** 

Jeroen van Wilgenburg

Shinichi Yamashita

Commissioning Editor

Edward Gordon

**Acquisition Editors** 

Joanne Fitzpatrick

**Content Development Editor** 

Shweta Pant

**Technical Editors** 

Indrajit A. Das

Pankaj Kadam

Copy Editors

Puja Lalwani

Alfida Paiva

Laxmi Subramanian

**Project Coordinator** 

Shipra Chawhan

**Proofreaders** 

**Bridget Braund** 

Maria Gould

Paul Hindle

Bernadette Watkins

Indexer

Priya Sane

**Production Coordinator** 

Nitesh Thakur

Cover Work

Nitesh Thakur

## **About the Author**

**Thilina Gunarathne** is a senior data scientist at KPMG LLP. He led the Hadoop-related efforts at Link Analytics before its acquisition by KPMG LLP. He has extensive experience in using Apache Hadoop and its related technologies for large-scale data-intensive computations. He coauthored the first edition of this book, *Hadoop MapReduce Cookbook*, with Dr. Srinath Perera.

Thilina has contributed to several open source projects at Apache Software Foundation as a member, committer, and a PMC member. He has also published many peer-reviewed research articles on how to extend the MapReduce model to perform efficient data mining and data analytics computations in the cloud. Thilina received his PhD and MSc degrees in computer science from Indiana University, Bloomington, USA, and received his bachelor of science degree in computer science and engineering from University of Moratuwa, Sri Lanka.

# **Acknowledgments**

I would like to thank my wife, Bimalee, my son, Kaveen, and my daughter, Yasali, for putting up with me for all the missing family time and for providing me with love and encouragement throughout the writing period. I would also like to thank my parents and siblings. Without their love, guidance, and encouragement, I would not be where I am today.

I really appreciate the contributions from my coauthor, Dr. Srinath Perera, for the first edition of this book. Many of his contributions from the first edition of this book have been adapted to the current book even though he wasn't able to coauthor this book due to his work and family commitments.

I would like to thank the Hadoop, HBase, Mahout, Pig, Hive, Sqoop, Nutch, and Lucene communities for developing great open source products. Thanks to Apache Software Foundation for fostering vibrant open source communities.

Big thanks to the editorial staff at Packt for providing me with the opportunity to write this book and feedback and guidance throughout the process. Thanks to the reviewers of this book for the many useful suggestions and corrections.

I would like to express my deepest gratitude to all the mentors I have had over the years, including Prof. Geoffrey Fox, Dr. Chris Groer, Dr. Sanjiva Weerawarana, Prof. Dennis Gannon, Prof. Judy Qiu, Prof. Beth Plale, and all my professors at Indiana University and University of Moratuwa for all the knowledge and guidance they gave me. Thanks to all my past and present colleagues for the many insightful discussions we've had and the knowledge they shared with me.

### **About the Author**

**Srinath Perera** (coauthor of the first edition of this book) is a senior software architect at WSO2 Inc., where he overlooks the overall WSO2 platform architecture with the CTO. He also serves as a research scientist at Lanka Software Foundation and teaches as a member of the visiting faculty at Department of Computer Science and Engineering, University of Moratuwa. He is a cofounder of Apache Axis2 open source project, and he has been involved with the Apache Web Service project since 2002 and is a member of Apache Software foundation and Apache Web Service project PMC. Srinath is also a committer of Apache open source projects Axis, Axis2, and Geronimo.

Srinath received his PhD and MSc in computer science from Indiana University, Bloomington, USA, and his bachelor of science in computer science and engineering from University of Moratuwa, Sri Lanka.

Srinath has authored many technical and peer-reviewed research articles; more details can be found on his website. He is also a frequent speaker at technical venues.

Srinath has worked with large-scale distributed systems for a long time. He closely works with big data technologies such as Hadoop and Cassandra daily. He also teaches a parallel programming graduate class at University of Moratuwa, which is primarily based on Hadoop.

I would like to thank my wife, Miyuru, and my parents, whose never-ending support keeps me going. I would also like to thank Sanjiva from WSO2 who encouraged us to make our mark even though project such as these are not in the job description. Finally, I would like to thank my colleagues at WSO2 for ideas and companionship that have shaped the book in many ways.

## **About the Reviewers**

**Skanda Bhargav** is an engineering graduate from Visvesvaraya Technological University (VTU), Belgaum, Karnataka, India. He did his majors in computer science engineering. He is currently employed with Happiest Minds Technologies, an MNC based out of Bangalore. He is a Cloudera-certified developer in Apache Hadoop. His interests are big data and Hadoop.

He has been a reviewer for the following books and a video, all by Packt Publishing:

- Instant MapReduce Patterns Hadoop Essentials How-to
- Hadoop Cluster Deployment
- Building Hadoop Clusters [Video]
- Cloudera Administration Handbook

I would like to thank my family for their immense support and faith in me throughout my learning stage. My friends have brought the confidence in me to a level that makes me bring out the best in myself. I am happy that God has blessed me with such wonderful people, without whom I wouldn't have tasted the success that I've achieved today.

**Randal Scott King** is a global consultant who specializes in big data and network architecture. His 15 years of experience in IT consulting has resulted in a client list that looks like a "Who's Who" of the Fortune 500. His recent projects include a complete network redesign for an aircraft manufacturer and an in-store video analytics pilot for a major home improvement retailer.

He lives with his children outside Atlanta, GA. You can visit his blog at www.randalscottking.com.

**Dmitry Spikhalskiy** currently holds the position of software engineer in a Russian social network service, Odnoklassniki, and is working on a search engine, video recommendation system, and movie content analysis.

Previously, Dmitry took part in developing Mind Labs' platform, infrastructure, and benchmarks for a high-load video conference and streaming service, which got "The biggest online-training in the world" Guinness world record with more than 12,000 participants. As a technical lead and architect, he launched a mobile social banking start-up called Instabank. He has also reviewed *Learning Google Guice* and *PostgreSQL 9 Admin Cookbook*, both by Packt Publishing.

Dmitry graduated from Moscow State University with an MSc degree in computer science, where he first got interested in parallel data processing, high-load systems, and databases.

**Jeroen van Wilgenburg** is a software craftsman at JPoint (http://www.jpoint.nl), a software agency based in the center of the Netherlands. Their main focus is on developing high-quality Java and Scala software with open source frameworks.

Currently, Jeroen is developing several big data applications with Hadoop, MapReduce, Storm, Spark, Kafka, MongoDB, and Elasticsearch.

Jeroen is a car enthusiast and likes to be outdoors, usually training for a triathlon. In his spare time, Jeroen writes about his work experience at http://vanwilgenburg.wordpress.com.

**Shinichi Yamashita** is a solutions architect at System Platform Sector in NTT DATA Corporation, Japan. He has more than 9 years of experience in software and middleware engineering (Apache, Tomcat, PostgreSQL, Hadoop Ecosystem, and Spark). Shinichi has written a few books on Hadoop in Japanese.

### www.PacktPub.com

#### Support files, eBooks, discount offers, and more

For support files and downloads related to your book, please visit www.PacktPub.com.

Did you know that Packt offers eBook versions of every book published, with PDF and ePub files available? You can upgrade to the eBook version at www.PacktPub.com and as a print book customer, you are entitled to a discount on the eBook copy. Get in touch with us at service@packtpub.com for more details.

At www.PacktPub.com, you can also read a collection of free technical articles, sign up for a range of free newsletters and receive exclusive discounts and offers on Packt books and eBooks.



https://www2.packtpub.com/books/subscription/packtlib

Do you need instant solutions to your IT questions? PacktLib is Packt's online digital book library. Here, you can search, access, and read Packt's entire library of books.

#### Why Subscribe?

- Fully searchable across every book published by Packt
- ▶ Copy and paste, print, and bookmark content
- On demand and accessible via a web browser

#### Free Access for Packt account holders

If you have an account with Packt at www.PacktPub.com, you can use this to access PacktLib today and view 9 entirely free books. Simply use your login credentials for immediate access.

# **Table of Contents**

Preface	1
Chapter 1: Getting Started with Hadoop v2	7
Introduction	7
Setting up Hadoop v2 on your local machine	14
Writing a WordCount MapReduce application, bundling it,	
and running it using the Hadoop local mode	15
Adding a combiner step to the WordCount MapReduce program	19
Setting up HDFS	21
Setting up Hadoop YARN in a distributed cluster environment	
using Hadoop v2	25
Setting up Hadoop ecosystem in a distributed cluster environment	
using a Hadoop distribution	28
HDFS command-line file operations	31
Running the WordCount program in a distributed cluster environment	32
Benchmarking HDFS using DFSIO	34
Benchmarking Hadoop MapReduce using TeraSort	35
Chapter 2: Cloud Deployments – Using Hadoop YARN on	
Cloud Environments	39
Introduction	39
Running Hadoop MapReduce v2 computations using Amazon	
Elastic MapReduce	41
Saving money using Amazon EC2 Spot Instances to execute EMR job flows	46
Executing a Pig script using EMR	48
Executing a Hive script using EMR	52
Creating an Amazon EMR job flow using the AWS Command Line Interface	55
Deploying an Apache HBase cluster on Amazon EC2 using EMR	59

Using EMR bootstrap actions to configure VMs for the Amazon EMR jobs Using Apache Whirr to deploy an Apache Hadoop cluster in a	63
cloud environment	65
Chapter 3: Hadoop Essentials – Configurations, Unit Tests,	
and Other APIs	69
Introduction	70
Optimizing Hadoop YARN and MapReduce configurations for	
cluster deployments	70
Shared user Hadoop clusters – using Fair and Capacity schedulers	74
Setting classpath precedence to user-provided JARs	76
Speculative execution of straggling tasks	77
Unit testing Hadoop MapReduce applications using MRUnit	78
Integration testing Hadoop MapReduce applications using	
MiniYarnCluster	82
Adding a new DataNode	85
Decommissioning DataNodes	87
Using multiple disks/volumes and limiting HDFS disk usage	88
Setting the HDFS block size	89
Setting the file replication factor	90
Using the HDFS Java API	92
Chapter 4: Developing Complex Hadoop MapReduce Applications	97
Introduction	98
Choosing appropriate Hadoop data types	99
Implementing a custom Hadoop Writable data type	102
Implementing a custom Hadoop key type	106
Emitting data of different value types from a Mapper	109
Choosing a suitable Hadoop InputFormat for your input data format	112
Adding support for new input data formats – implementing	
a custom InputFormat	114
Formatting the results of MapReduce computations – using	
Hadoop OutputFormats	118
Writing multiple outputs from a MapReduce computation	120
Hadoop intermediate data partitioning	123
Secondary sorting – sorting Reduce input values	125
Broadcasting and distributing shared resources to tasks in a	
MapReduce job – Hadoop DistributedCache	129
Using Hadoop with legacy applications – Hadoop streaming	132
Adding dependencies between MapReduce jobs	135
Hadoop counters to report custom metrics	137
manach acquirere to take to another morning	201

Chapter 5: Analytics	141
Introduction	141
Simple analytics using MapReduce	142
Performing GROUP BY using MapReduce	145
Calculating frequency distributions and sorting using MapReduce	148
Plotting the Hadoop MapReduce results using gnuplot	150
Calculating histograms using MapReduce	152
Calculating Scatter plots using MapReduce	155
Parsing a complex dataset with Hadoop	158
Joining two datasets using MapReduce	163
Chapter 6: Hadoop Ecosystem - Apache Hive	169
Introduction	170
Getting started with Apache Hive	171
Creating databases and tables using Hive CLI	172
Simple SQL-style data querying using Apache Hive	177
Creating and populating Hive tables and views using Hive query results	182
Utilizing different storage formats in Hive - storing table data	
using ORC files	183
Using Hive built-in functions	185
Hive batch mode - using a query file	187
Performing a join with Hive	189
Creating partitioned Hive tables	191
Writing Hive User-defined Functions (UDF)	193
HCatalog – performing Java MapReduce computations on	200
data mapped to Hive tables	195
HCatalog – writing data to Hive tables from Java	100
MapReduce computations	198
Chapter 7: Hadoop Ecosystem II – Pig, HBase, Mahout, and Sqoop	201
Introduction	201
Getting started with Apache Pig	203
Joining two datasets using Pig	206
Accessing a Hive table data in Pig using HCatalog	208
Getting started with Apache HBase	210
Data random access using Java client APIs	213
Running MapReduce jobs on HBase	214
Using Hive to insert data into HBase tables	215
Getting started with Apache Mahout	218
Running K-means with Mahout	219
Importing data to HDFS from a relational database using Apache Sqoop	220
Exporting data from HDFS to a relational database using Apache Sqoop	225

Chapter 8: Searching and Indexing	229
Introduction	229
Generating an inverted index using Hadoop MapReduce	230
Intradomain web crawling using Apache Nutch	235
Indexing and searching web documents using Apache Solr	240
Configuring Apache HBase as the backend data store for Apache Nutch	243
Whole web crawling with Apache Nutch using a Hadoop/HBase cluster	246
Elasticsearch for indexing and searching	249
Generating the in-links graph for crawled web pages	251
Chapter 9: Classifications, Recommendations,	
and Finding Relationships	255
Introduction	255
Performing content-based recommendations	256
Classification using the naïve Bayes classifier	261
Assigning advertisements to keywords using the Adwords	
balance algorithm	265
Chapter 10: Mass Text Data Processing	273
Introduction	273
Data preprocessing using Hadoop streaming and Python	274
De-duplicating data using Hadoop streaming	277
Loading large datasets to an Apache HBase data store - importtsv	
and bulkload	279
Creating TF and TF-IDF vectors for the text data	283
Clustering text data using Apache Mahout	288
Topic discovery using Latent Dirichlet Allocation (LDA)	290
Document classification using Mahout Naive Bayes Classifier	293
Index	297

# **Preface**

We are currently facing an avalanche of data, and this data contains many insights that hold the keys to success or failure in the data-driven world. Next generation Hadoop (v2) offers a cutting-edge platform to store and analyze these massive data sets and improve upon the widely used and highly successful Hadoop MapReduce v1. The recipes that will help you analyze large and complex datasets with next generation Hadoop MapReduce will provide you with the skills and knowledge needed to process large and complex datasets using the next generation Hadoop ecosystem.

This book presents many exciting topics such as MapReduce patterns using Hadoop to solve analytics, classifications, and data indexing and searching. You will also be introduced to several Hadoop ecosystem components including Hive, Pig, HBase, Mahout, Nutch, and Sqoop.

This book introduces you to simple examples and then dives deep to solve in-depth big data use cases. This book presents more than 90 ready-to-use Hadoop MapReduce recipes in a simple and straightforward manner, with step-by-step instructions and real-world examples.

### What this book covers

Chapter 1, Getting Started with Hadoop v2, introduces Hadoop MapReduce, YARN, and HDFS, and walks through the installation of Hadoop v2.

Chapter 2, Cloud Deployments – Using Hadoop Yarn on Cloud Environments, explains how to use Amazon Elastic MapReduce (EMR) and Apache Whirr to deploy and execute Hadoop MapReduce, Pig, Hive, and HBase computations on cloud infrastructures.

Chapter 3, Hadoop Essentials – Configurations, Unit Tests, and Other APIs, introduces basic Hadoop YARN and HDFS configurations, HDFS Java API, and unit testing methods for MapReduce applications.

Chapter 4, Developing Complex Hadoop MapReduce Applications, introduces you to several advanced Hadoop MapReduce features that will help you develop highly customized and efficient MapReduce applications.

ra		

Chapter 5, Analytics, explains how to perform basic data analytic operations using Hadoop MapReduce.

Chapter 6, Hadoop Ecosystem – Apache Hive, introduces Apache Hive, which provides data warehouse capabilities on top of Hadoop, using a SQL-like query language.

Chapter 7, Hadoop Ecosystem II – Pig, HBase, Mahout, and Sqoop, introduces the Apache Pig data flow style data-processing language, Apache HBase NoSQL data storage, Apache Mahout machine learning and data-mining toolkit, and Apache Sqoop bulk data transfer utility to transfer data between Hadoop and the relational databases.

Chapter 8, Searching and Indexing, introduces several tools and techniques that you can use with Apache Hadoop to perform large-scale searching and indexing.

Chapter 9, Classifications, Recommendations, and Finding Relationships, explains how to implement complex algorithms such as classifications, recommendations, and finding relationships using Hadoop.

Chapter 10, Mass Text Data Processing, explains how to use Hadoop and Mahout to process large text datasets and how to perform data preprocessing and loading of operations using Hadoop.

### What you need for this book

You need a moderate knowledge of Java and access to the Internet and a computer that runs a Linux operating system.

#### Who this book is for

If you are a big data enthusiast and wish to use Hadoop v2 to solve your problems, then this book is for you. This book is for Java programmers with little to moderate knowledge of Hadoop MapReduce. This is also a one-stop reference for developers and system admins who want to quickly get up to speed with using Hadoop v2. It would be helpful to have a basic knowledge of software development using Java and a basic working knowledge of Linux.

#### Conventions

In this book, you will find a number of styles of text that distinguish between different kinds of information. Here are some examples of these styles, and an explanation of their meaning.

Code words in text, database table names, folder names, file are shown as follows: "The following are the descriptions of the properties we used in the hadoop.properties file."

A block of code is set as follows:

```
Path file = new Path("demo.txt");
FSDataOutputStream outStream = fs.create(file);
outStream.writeUTF("Welcome to HDFS Java API!!!");
outStream.close();
```

When we wish to draw your attention to a particular part of a code block, the relevant lines or items are set in bold:

```
Job job = Job.getInstance(getConf(), "MLReceiveReplyProcessor");
job.setJarByClass(CountReceivedRepliesMapReduce.class);
job.setMapperClass(AMapper.class);
job.setReducerClass(AReducer.class);
job.setNumReduceTasks(numReduce);

job.setOutputKeyClass(Text.class);
job.setOutputValueClass(Text.class);
job.setInputFormatClass(MBoxFileInputFormat.class);
FileInputFormat.setInputPaths(job, new Path(inputPath));
FileOutputFormat.setOutputPath(job, new Path(outputPath));
int exitStatus = job.waitForCompletion(true) ? 0 : 1;
```

Any command-line input or output is written as follows:

```
205.212.115.106 - - [01/Jul/1995:00:00:12 -0400] "GET /shuttle/countdown/countdown.html HTTP/1.0" 200 3985
```

**New terms** and **important words** are shown in bold. Words that you see on the screen, in menus or dialog boxes for example, appear in the text like this: "Select **Custom Action** in the **Add Bootstrap Actions** drop-down box. Click on **Configure and add.**"



Warnings or important notes appear in a box like this.



Tips and tricks appear like this.

3