



Quick answers to common problems

# Hadoop MapReduce Cookbook

Recipes for analyzing large and complex datasets with Hadoop MapReduce

**Srinath Perera**  
**Thilina Gunarathne**

**[PACKT]** open source\*  
PUBLISHING community experience distilled

---

# Hadoop MapReduce Cookbook

Recipes for analyzing large and complex datasets with Hadoop MapReduce

**Srinath Perera**

**Thilina Gunarathne**

**[PACKT]** open source   
PUBLISHING community experience distilled

BIRMINGHAM - MUMBAI

---

## **Hadoop MapReduce Cookbook**

Copyright © 2013 Packt Publishing

All rights reserved. No part of this book may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, without the prior written permission of the publisher, except in the case of brief quotations embedded in critical articles or reviews.

Every effort has been made in the preparation of this book to ensure the accuracy of the information presented. However, the information contained in this book is sold without warranty, either express or implied. Neither the authors, nor Packt Publishing, and its dealers and distributors will be held liable for any damages caused or alleged to be caused directly or indirectly by this book.

Packt Publishing has endeavored to provide trademark information about all of the companies and products mentioned in this book by the appropriate use of capitals. However, Packt Publishing cannot guarantee the accuracy of this information.

First published: February 2013

Production Reference: 2250113

Published by Packt Publishing Ltd.  
Livery Place  
35 Livery Street  
Birmingham B3 2PB, UK.

ISBN 978-1-84951-728-7

[www.packtpub.com](http://www.packtpub.com)

Cover Image by J.Blaminsky (milak6@wp.pl)

---

# Credits

**Authors**

Srinath Perera  
Thilina Gunarathne

**Reviewers**

Masatake Iwasaki  
Shinichi Yamashita

**Acquisition Editor**

Robin de Jongh

**Lead Technical Editor**

Arun Nadar

**Technical Editors**

Vrinda Amberkar  
Dennis John  
Dominic Pereira

**Project Coordinator**

Amey Sawant

**Proofreader**

Mario Cecere

**Indexer**

Hemangini Bari

**Graphics**

Valentina D'Silva

**Production Coordinator**

Arvindkumar Gupta

**Cover Work**

Arvindkumar Gupta

---

# About the Authors

**Srinath Perera** 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 visiting faculty at Department of Computer Science and Engineering, University of Moratuwa. He is a co-founder 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.

He received his Ph.D. and M.Sc. in Computer Sciences from Indiana University, Bloomington, USA and received 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, and more detail can be found from his website. He is also a frequent speaker at technical venues.

He 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 also like to thanks Sanjiva from WSO2 who encourage us to make our mark even though project like these are not in the job description. Finally I would like to thank my colleges at WSO2 for ideas and companionship that have shaped the book in many ways.

---

---

**Thilina Gunarathne** is a Ph.D. candidate at the School of Informatics and Computing of Indiana University. He has extensive experience in using Apache Hadoop and related technologies for large-scale data intensive computations. His current work focuses on developing technologies to perform scalable and efficient large-scale data intensive computations on cloud environments.

Thilina has published many articles and peer reviewed research papers in the areas of distributed and parallel computing, including several papers on extending MapReduce model to perform efficient data mining and data analytics computations on clouds. Thilina is a regular presenter in both academic as well as industry settings.

Thilina has contributed to several open source projects at Apache Software Foundation as a committer and a PMC member since 2005. Before starting the graduate studies, Thilina worked as a Senior Software Engineer at WSO2 Inc., focusing on open source middleware development. Thilina received his B.Sc. in Computer Science and Engineering from University of Moratuwa, Sri Lanka, in 2006 and received his M.Sc. in Computer Science from Indiana University, Bloomington, in 2009. Thilina expects to receive his doctorate in the field of distributed and parallel computing in 2013.

---

This book would not have been a success without the direct and indirect help from many people. Thanks to my wife and my son for putting up with me for all the missing family times and for providing me with love and encouragement throughout the writing period. Thanks to my parents, without whose love, guidance and encouragement, I would not be where I am today.

Thanks to my advisor Prof. Geoffrey Fox for his excellent guidance and providing me with the environment to work on Hadoop and related technologies. Thanks to the HBase, Mahout, Pig, Hive, Nutch, and Lucene communities for developing great open source products. Thanks to Apache Software Foundation for fostering vibrant open source communities.

Thanks to the editorial staff at Packt, for providing me the opportunity to write this book and for providing feedback and guidance throughout the process. Thanks to the reviewers for reviewing this book, catching my mistakes, and for the many useful suggestions.

Thanks to all of my past and present mentors and teachers, including Dr. Sanjiva Weerawarana of WSO2, Prof. Dennis Gannon, Prof. Judy Qiu, Prof. Beth Plale, 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 many insightful discussions and the knowledge they shared with me.

---

---

# About the Reviewers

**Masatake Iwasaki** is Software Engineer at NTT DATA Corporation. He provides technical consultation for Open Source software such as Hadoop, HBase, and PostgreSQL.

**Shinichi Yamashita** is a Chief Engineer at OSS professional service unit in NTT DATA Corporation in Japan. He has more than seven years' experience in software and middleware (Apache, Tomcat, PostgreSQL, and Hadoop eco system) engineering. NTT DATA is your Innovation Partner anywhere around the world. It provides professional services from consulting, and system development to business IT outsourcing. In Japan, he has authored some books on Hadoop.

---

I thank my co-workers.

---

---

# www.PacktPub.com

## Support files, eBooks, discount offers and more

You might want to visit [www.PacktPub.com](http://www.PacktPub.com) for support files and downloads related to your book.

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](http://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](mailto:service@packtpub.com) for more details.

At [www.PacktPub.com](http://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.



<http://PacktLib.PacktPub.com>

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

## Free Access for Packt account holders

If you have an account with Packt at [www.PacktPub.com](http://www.PacktPub.com), you can use this to access PacktLib today and view nine entirely free books. Simply use your login credentials for immediate access.



---

# Table of Contents

<b>Preface</b>	<b>1</b>
<b>Chapter 1: Getting Hadoop Up and Running in a Cluster</b>	<b>5</b>
Introduction	5
Setting up Hadoop on your machine	6
Writing a WordCount MapReduce sample, bundling it, and running it using standalone Hadoop	7
Adding the combiner step to the WordCount MapReduce program	12
Setting up HDFS	13
Using HDFS monitoring UI	17
HDFS basic command-line file operations	18
Setting Hadoop in a distributed cluster environment	20
Running the WordCount program in a distributed cluster environment	24
Using MapReduce monitoring UI	26
<b>Chapter 2: Advanced HDFS</b>	<b>29</b>
Introduction	29
Benchmarking HDFS	30
Adding a new DataNode	31
Decommissioning DataNodes	33
Using multiple disks/volumes and limiting HDFS disk usage	34
Setting HDFS block size	35
Setting the file replication factor	36
Using HDFS Java API	38
Using HDFS C API (libhdfs)	42
Mounting HDFS (Fuse-DFS)	46
Merging files in HDFS	49

<b>Chapter 3: Advanced Hadoop MapReduce Administration</b>	<b>51</b>
Introduction	51
Tuning Hadoop configurations for cluster deployments	52
Running benchmarks to verify the Hadoop installation	54
Reusing Java VMs to improve the performance	56
Fault tolerance and speculative execution	56
Debug scripts – analyzing task failures	57
Setting failure percentages and skipping bad records	60
Shared-user Hadoop clusters – using fair and other schedulers	62
Hadoop security – integrating with Kerberos	63
Using the Hadoop Tool interface	69
<b>Chapter 4: Developing Complex Hadoop MapReduce Applications</b>	<b>73</b>
Introduction	74
Choosing appropriate Hadoop data types	74
Implementing a custom Hadoop Writable data type	77
Implementing a custom Hadoop key type	80
Emitting data of different value types from a mapper	83
Choosing a suitable Hadoop InputFormat for your input data format	87
Adding support for new input data formats – implementing a custom InputFormat	90
Formatting the results of MapReduce computations – using Hadoop OutputFormats	93
Hadoop intermediate (map to reduce) data partitioning	95
Broadcasting and distributing shared resources to tasks in a MapReduce job – Hadoop DistributedCache	97
Using Hadoop with legacy applications – Hadoop Streaming	101
Adding dependencies between MapReduce jobs	104
Hadoop counters for reporting custom metrics	106
<b>Chapter 5: Hadoop Ecosystem</b>	<b>109</b>
Introduction	109
Installing HBase	110
Data random access using Java client APIs	113
Running MapReduce jobs on HBase (table input/output)	115
Installing Pig	118
Running your first Pig command	119
Set operations (join, union) and sorting with Pig	121
Installing Hive	123
Running a SQL-style query with Hive	124
Performing a join with Hive	127
Installing Mahout	129

Running K-means with Mahout	130
Visualizing K-means results	132
<b>Chapter 6: Analytics</b>	<b>135</b>
Introduction	135
Simple analytics using MapReduce	136
Performing Group-By using MapReduce	140
Calculating frequency distributions and sorting using MapReduce	143
Plotting the Hadoop results using GNU Plot	145
Calculating histograms using MapReduce	147
Calculating scatter plots using MapReduce	151
Parsing a complex dataset with Hadoop	154
Joining two datasets using MapReduce	159
<b>Chapter 7: Searching and Indexing</b>	<b>165</b>
Introduction	165
Generating an inverted index using Hadoop MapReduce	166
Intra-domain web crawling using Apache Nutch	170
Indexing and searching web documents using Apache Solr	174
Configuring Apache HBase as the backend data store for Apache Nutch	177
Deploying Apache HBase on a Hadoop cluster	180
Whole web crawling with Apache Nutch using a Hadoop/HBase cluster	182
ElasticSearch for indexing and searching	185
Generating the in-links graph for crawled web pages	187
<b>Chapter 8: Classifications, Recommendations, and Finding Relationships</b>	<b>191</b>
Introduction	191
Content-based recommendations	192
Hierarchical clustering	198
Clustering an Amazon sales dataset	201
Collaborative filtering-based recommendations	205
Classification using Naive Bayes Classifier	208
Assigning advertisements to keywords using the Adwords balance algorithm	214
<b>Chapter 9: Mass Text Data Processing</b>	<b>223</b>
Introduction	223
Data preprocessing (extract, clean, and format conversion) using Hadoop Streaming and Python	224
Data de-duplication using Hadoop Streaming	227
Loading large datasets to an Apache HBase data store using importtsv and bulkload tools	229

<b>Creating TF and TF-IDF vectors for the text data</b>	<b>234</b>
<b>Clustering the text data</b>	<b>238</b>
<b>Topic discovery using Latent Dirichlet Allocation (LDA)</b>	<b>241</b>
<b>Document classification using Mahout Naive Bayes classifier</b>	<b>244</b>
<b>Chapter 10: Cloud Deployments: Using Hadoop on Clouds</b>	<b>247</b>
<b>Introduction</b>	<b>247</b>
<b>Running Hadoop MapReduce computations using Amazon Elastic MapReduce (EMR)</b>	<b>248</b>
<b>Saving money by using Amazon EC2 Spot Instances to execute EMR job flows</b>	<b>252</b>
<b>Executing a Pig script using EMR</b>	<b>253</b>
<b>Executing a Hive script using EMR</b>	<b>256</b>
<b>Creating an Amazon EMR job flow using the Command Line Interface</b>	<b>260</b>
<b>Deploying an Apache HBase Cluster on Amazon EC2 cloud using EMR</b>	<b>263</b>
<b>Using EMR Bootstrap actions to configure VMs for the Amazon EMR jobs</b>	<b>268</b>
<b>Using Apache Whirr to deploy an Apache Hadoop cluster in a cloud environment</b>	<b>270</b>
<b>Using Apache Whirr to deploy an Apache HBase cluster in a cloud environment</b>	<b>274</b>
<b>Index</b>	<b>277</b>

---

---

# Preface

*Hadoop MapReduce Cookbook* helps readers learn to process large and complex datasets. The book starts in a simple manner, but still provides in-depth knowledge of Hadoop. It is a simple one-stop guide on how to get things done. It has 90 recipes, presented in a simple and straightforward manner, with step-by-step instructions and real world examples.

This product includes software developed at The Apache Software Foundation (<http://www.apache.org/>).

## What this book covers

*Chapter 1, Getting Hadoop Up and Running in a Cluster*, explains how to install and run Hadoop both as a single node as well as a cluster.

*Chapter 2, Advanced HDFS*, introduces a set of advanced HDFS operations that would be useful when performing large-scale data processing with Hadoop MapReduce as well as with non-MapReduce use cases.

*Chapter 3, Advanced Hadoop MapReduce Administration*, explains how to change configurations and security of a Hadoop installation and how to debug.

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

*Chapter 5, Hadoop Ecosystem*, introduces the other projects related to Hadoop such as HBase, Hive, and Pig.

*Chapter 6, Analytics*, explains how to calculate basic analytics using Hadoop.

*Chapter 7, Searching and Indexing*, introduces you to several tools and techniques that you can use with Apache Hadoop to perform large-scale searching and indexing.

*Chapter 8, Classifications, Recommendations, and Finding Relationships*, explains how to implement complex algorithms such as classifications, recommendations, and finding relationships using Hadoop.

*Chapter 9, Mass Text Data Processing*, explains how to use Hadoop and Mahout to process large text datasets, and how to perform data preprocessing and loading operations using Hadoop.

*Chapter 10, Cloud Deployments: Using Hadoop on Clouds*, 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.

## What you need for this book

All you need is access to a computer running Linux Operating system, and Internet. Also, Java knowledge is required.

## Who this book is for

For big data enthusiasts and would be Hadoop programmers. The books for Java programmers who either have not worked with Hadoop at all, or who knows Hadoop and MapReduce but want to try out things and get into details. It is also a one-stop reference for most of your Hadoop tasks.

## 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 are shown as follows: "From this point onward, we shall call the unpacked Hadoop directory `HADOOP_HOME`."

A block of code is set as follows:

```
public void map(Object key, Text value, Context context)
    throws IOException, InterruptedException
{
    StringTokenizer itr = new StringTokenizer(value.toString());
    while (itr.hasMoreTokens())
    {
        word.set(itr.nextToken());
        context.write(word, new IntWritable(1));
    }
}
```

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

```
>tar -zxvf hadoop-1.x.x.tar.gz
```

**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: "Create a S3 bucket to upload the input data by clicking on **Create Bucket**".

 Warnings or important notes appear in a box like this. 

 Tips and tricks appear like this. 

## Reader feedback

Feedback from our readers is always welcome. Let us know what you think about this book—what you liked or may have disliked. Reader feedback is important for us to develop titles that you really get the most out of.

To send us general feedback, simply send an e-mail to [feedback@packtpub.com](mailto:feedback@packtpub.com), and mention the book title via the subject of your message.

If there is a topic that you have expertise in and you are interested in either writing or contributing to a book, see our author guide on [www.packtpub.com/authors](http://www.packtpub.com/authors).

## Customer support

Now that you are the proud owner of a Packt book, we have a number of things to help you to get the most from your purchase.

## Downloading the example code

You can download the example code files for all Packt books you have purchased from your account at <http://www.PacktPub.com>. If you purchased this book elsewhere, you can visit <http://www.PacktPub.com/support> and register to have the files e-mailed directly to you.

## Errata

Although we have taken every care to ensure the accuracy of our content, mistakes do happen. If you find a mistake in one of our books—maybe a mistake in the text or the code—we would be grateful if you would report this to us. By doing so, you can save other readers from frustration and help us improve subsequent versions of this book. If you find any errata, please report them by visiting <http://www.packtpub.com/support>, selecting your book, clicking on the **errata submission form** link, and entering the details of your errata. Once your errata are verified, your submission will be accepted and the errata will be uploaded on our website, or added to any list of existing errata, under the Errata section of that title. Any existing errata can be viewed by selecting your title from <http://www.packtpub.com/support>.

## Piracy

Piracy of copyright material on the Internet is an ongoing problem across all media. At Packt, we take the protection of our copyright and licenses very seriously. If you come across any illegal copies of our works, in any form, on the Internet, please provide us with the location address or website name immediately so that we can pursue a remedy.

Please contact us at [copyright@packtpub.com](mailto:copyright@packtpub.com) with a link to the suspected pirated material.

We appreciate your help in protecting our authors, and our ability to bring you valuable content.

## Questions

You can contact us at [questions@packtpub.com](mailto:questions@packtpub.com) if you are having a problem with any aspect of the book, and we will do our best to address it.

---

# 1

## Getting Hadoop Up and Running in a Cluster

In this chapter, we will cover:

- ▶ Setting up Hadoop on your machine
- ▶ Writing the WordCount MapReduce sample, bundling it, and running it using standalone Hadoop
- ▶ Adding the combiner step to the WordCount MapReduce program
- ▶ Setting up HDFS
- ▶ Using the HDFS monitoring UI
- ▶ HDFS basic command-line file operations
- ▶ Setting Hadoop in a distributed cluster environment
- ▶ Running the WordCount program in a distributed cluster environment
- ▶ Using the MapReduce monitoring UI

### Introduction

For many years, users who want to store and analyze data would store the data in a database and process it via SQL queries. The Web has changed most of the assumptions of this era. On the Web, the data is unstructured and large, and the databases can neither capture the data into a schema nor scale it to store and process it.

Google was one of the first organizations to face the problem, where they wanted to download the whole of the Internet and index it to support search queries. They built a framework for large-scale data processing borrowing from the "map" and "reduce" functions of the functional programming paradigm. They called the paradigm **MapReduce**.

Hadoop is the most widely known and widely used implementation of the MapReduce paradigm. This chapter introduces Hadoop, describes how to install Hadoop, and shows you how to run your first MapReduce job with Hadoop.

Hadoop installation consists of four types of nodes—a **NameNode**, **DataNodes**, a **JobTracker**, and **TaskTracker** HDFS nodes (NameNode and DataNodes) provide a distributed filesystem where the JobTracker manages the jobs and TaskTrackers run tasks that perform parts of the job. Users submit MapReduce jobs to the JobTracker, which runs each of the Map and Reduce parts of the initial job in TaskTrackers, collects results, and finally emits the results.

Hadoop provides three installation choices:

- ▶ **Local mode:** This is an unzip and run mode to get you started right away where all parts of Hadoop run within the same JVM
- ▶ **Pseudo distributed mode:** This mode will be run on different parts of Hadoop as different Java processors, but within a single machine
- ▶ **Distributed mode:** This is the real setup that spans multiple machines

We will discuss the local mode in the first three recipes, and Pseudo distributed and distributed modes in the last three recipes.

## Setting up Hadoop on your machine

This recipe describes how to run Hadoop in the local mode.

### Getting ready

Download and install Java 1.6 or higher version from <http://www.oracle.com/technetwork/java/javase/downloads/index.html>.

### How to do it...

Now let us do the Hadoop installation:

1. Download the most recent Hadoop 1.0 branch distribution from <http://hadoop.apache.org/>.
2. Unzip the Hadoop distribution using the following command. You will have to change the `x.x` in the filename with the actual release you have downloaded. If you are using Windows, you should use your favorite archive program such as WinZip or WinRAR for extracting the distribution. From this point onward, we shall call the unpacked Hadoop directory `HADOOP_HOME`.

```
>tar -zxvf hadoop-1.x.x.tar.gz
```

- You can use Hadoop local mode after unzipping the distribution. Your installation is done. Now, you can run Hadoop jobs through `bin/hadoop` command, and we will elaborate that further in the next recipe.

### How it works...

Hadoop local mode does not start any servers but does all the work within the same JVM. When you submit a job to Hadoop in the local mode, that job starts a JVM to run the job, and that JVM carries out the job. The output and the behavior of the job is the same as a distributed Hadoop job, except for the fact that the job can only use the current node for running tasks. In the next recipe, we will discover how to run a MapReduce program using the unzipped Hadoop distribution.

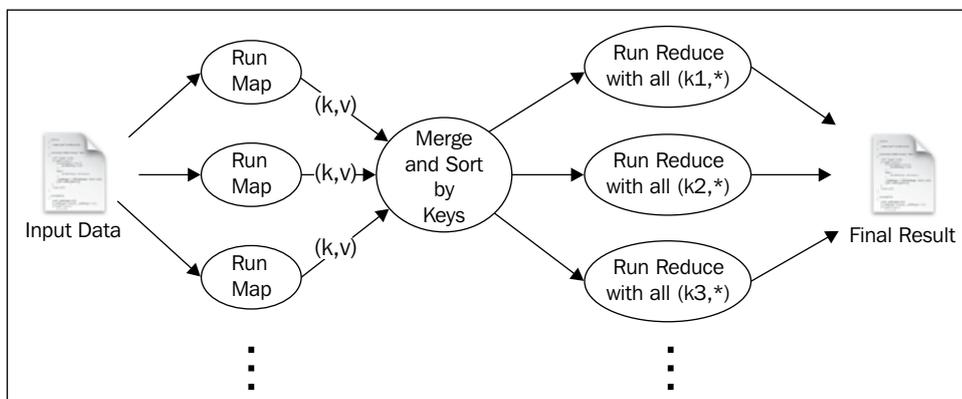


#### Downloading the example code

You can download the example code files for all Packt books you have purchased from your account at <http://www.packtpub.com>. If you purchased this book elsewhere, you can visit <http://www.packtpub.com/support> and register to have the files e-mailed directly to you.

## Writing a WordCount MapReduce sample, bundling it, and running it using standalone Hadoop

This recipe explains how to write a simple MapReduce program and how to execute it.



To run a MapReduce job, users should furnish a `map` function, a `reduce` function, input data, and an output data location. When executed, Hadoop carries out the following steps:

1. Hadoop breaks the input data into multiple data items by new lines and runs the `map` function once for each data item, giving the item as the input for the function. When executed, the `map` function outputs one or more key-value pairs.
2. Hadoop collects all the key-value pairs generated from the `map` function, sorts them by the key, and groups together the values with the same key.
3. For each distinct key, Hadoop runs the `reduce` function once while passing the key and list of values for that key as input.
4. The `reduce` function may output one or more key-value pairs, and Hadoop writes them to a file as the final result.

## Getting ready

From the source code available with this book, select the source code for the first chapter, `chapter1_src.zip`. Then, set it up with your favorite **Java Integrated Development Environment (IDE)**; for example, Eclipse. You need to add the `hadoop-core` JAR file in `HADOOP_HOME` and all other JAR files in the `HADOOP_HOME/lib` directory to the classpath of the IDE.

Download and install Apache Ant from <http://ant.apache.org/>.

## How to do it...

Now let us write our first Hadoop MapReduce program.

1. The `WordCount` sample uses MapReduce to count the number of word occurrences within a set of input documents. Locate the sample code from `src/chapter1/Wordcount.java`. The code has three parts—mapper, reducer, and the main program.
2. The mapper extends from the `org.apache.hadoop.mapreduce.Mapper` interface. When Hadoop runs, it receives each new line in the input files as an input to the mapper. The `map` function breaks each line into substrings using whitespace characters such as the separator, and for each token (word) emits `(word, 1)` as the output.

```
public void map(Object key, Text value, Context context
                ) throws IOException, InterruptedException
{
    StringTokenizer itr = new StringTokenizer(value.toString());
    while (itr.hasMoreTokens())
```

```

    {
        word.set(itr.nextToken());
        context.write(word, new IntWritable(1));
    }
}

```

3. The reduce function receives all the values that have the same key as the input, and it outputs the key and the number of occurrences of the key as the output.

```

public void reduce(Text key, Iterable<IntWritable> values,
                  Context context
                  ) throws IOException, InterruptedException
{
    int sum = 0;
    for (IntWritable val : values)
    {
        sum += val.get();
    }
    result.set(sum);
    context.write(key, result);
}

```

4. The main program puts the configuration together and submits the job to Hadoop.

```

Configuration conf = new Configuration();
String[] otherArgs = new GenericOptionsParser(conf, args).
getRemainingArgs();
if (otherArgs.length != 2) {
    System.err.println("Usage: wordcount <in><out>");
    System.exit(2);
}
Job job = new Job(conf, "word count");
job.setJarByClass(WordCount.class);
job.setMapperClass(TokenMapper.class);
//Uncomment this to
//job.setCombinerClass(IntSumReducer.class);
job.setReducerClass(IntSumReducer.class);
job.setOutputKeyClass(Text.class);
job.setOutputValueClass(IntWritable.class);
FileInputFormat.addInputPath(job, new Path(otherArgs[0]));
FileOutputFormat.setOutputPath(job, new Path(otherArgs[1]));
System.exit(job.waitForCompletion(true) ? 0 : 1);

```

5. You can compile the sample by running the following command, which uses Apache Ant, from the root directory of the sample code:

```
>ant build
```

If you have not done this already, you should install Apache Ant by following the instructions given at <http://ant.apache.org/manual/install.html>. Alternatively, you can use the compiled JAR file included with the source code.

6. Change the directory to `HADOOP_HOME`, and copy the `hadoop-cookbook-chapter1.jar` file to the `HADOOP_HOME` directory. To be used as the input, create a directory called `input` under `HADOOP_HOME` and copy the `README.txt` file to the directory. Alternatively, you can copy any text file to the `input` directory.
7. Run the sample using the following command. Here, `chapter1.WordCount` is the name of the `main` class we need to run. When you have run the command, you will see the following terminal output:

```
>bin/hadoop jar hadoop-cookbook-chapter1.jar chapter1.WordCount
input output

12/04/11 08:12:44 INFO input.FileInputFormat: Total input paths to
process : 16

12/04/11 08:12:45 INFO mapred.JobClient: Running job: job_
local_0001

12/04/11 08:12:45 INFO mapred.Task: Task:attempt_
local_0001_m_000000_0 is done. And is in the process of committing
.....

.....

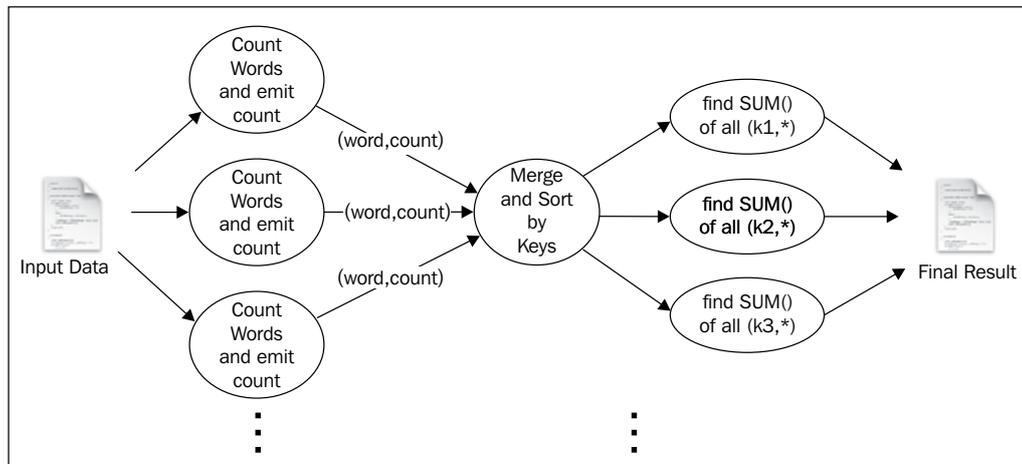
12/04/11 08:13:37 INFO mapred.JobClient: Job complete: job_
local_0001

.....
```

8. The output directory will have a file named like `part-r-XXXXX`, which will have the count of each word in the document. Congratulations! You have successfully run your first MapReduce program.

### How it works...

In the preceding sample, MapReduce worked in the local mode without starting any servers and using the local filesystem as the storage system for inputs, outputs, and working data. The following diagram shows what happened in the WordCount program under the covers:



The workflow is as follows:

1. Hadoop reads the input, breaks it by new line characters as the separator and then runs the map function passing each line as an argument.
2. The map function tokenizes the line, and for each token (word), emits a key value pair (word, 1).
3. Hadoop collects all the (word, 1) pairs, sorts them by the word, groups all the values emitted against each unique key, and invokes the reduce once for each unique key passing the key and values for that key as an argument.
4. The reduce function counts the number of occurrences of each word using the values and emits it as a key-value pair.
5. Hadoop writes the final output to the output directory.

### There's more...

As an optional step, copy the `input` directory to the top level of the IDE-based project (Eclipse project) that you created for samples. Now you can run the `wordCount` class directly from your IDE passing `input` `output` as arguments. This will run the sample the same as before. Running MapReduce jobs from IDE in this manner is very useful for debugging your MapReduce jobs.

Although you ran the sample with Hadoop installed in your local machine, you can run it using distributed Hadoop cluster setup with a HDFS-distributed filesystem. The recipes of this chapter, *Setting up HDFS* and *Setting Hadoop in a distributed cluster environment* will discuss how to run this sample in a distributed setup.

- [download online His Bright Light: The Story of Nick Traina for free](#)
- [The Eudaemonic Pie book](#)
- [read online Through the Eyes of Tiger Cubs: Views of Asia's Next Generation](#)
- [Battle for Atlantis \(Atlantis, Book 6\) pdf](#)
- [click The Story of Philosophy: The Essential Guide to the History of Western Philosophy.pdf, azw \(kindle\)](#)
- [read online Evolution 2.0: Breaking the Deadlock Between Darwin and Design](#)
  
- <http://sidenoter.com/?ebooks/Soldier-of-the-Mist--Soldier-of-the-Mist--Book-1-.pdf>
- <http://interactmg.com/ebooks/The-Biological-Evolution-of-Religious-Mind-and-Behavior--The-Frontiers-Collection-.pdf>
- <http://diy-chirol.com/lib/Napalm--An-American-Biography.pdf>
- <http://creativebeard.ru/freebooks/Mood-Indigo.pdf>
- <http://crackingscience.org/?library/The-Story-of-Philosophy--The-Essential-Guide-to-the-History-of-Western-Philosophy.pdf>
- <http://omarnajmi.com/library/Evolution-2-0--Breaking-the-Deadlock-Between-Darwin-and-Design.pdf>