



# Beginning Apache Pig

Big Data Processing Made Easy

—  
Balaswamy Vaddeman

Apress®

# Beginning Apache Pig

Big Data Processing Made Easy



**Balaswamy Vaddeman**

Apress®

## ***Beginning Apache Pig: Big Data Processing Made Easy***

Balaswamy Vaddeman  
Hyderabad, Andhra Pradesh, India

ISBN-13 (pbk): 978-1-4842-2336-9

ISBN-13 (electronic): 978-1-4842-2337-6

DOI 10.1007/978-1-4842-2337-6

Library of Congress Control Number: 2016961514

Copyright © 2016 by Balaswamy Vaddeman

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

Trademarked names, logos, and images may appear in this book. Rather than use a trademark symbol with every occurrence of a trademarked name, logo, or image we use the names, logos, and images only in an editorial fashion and to the benefit of the trademark owner, with no intention of infringement of the trademark.

The use in this publication of trade names, trademarks, service marks, and similar terms, even if they are not identified as such, is not to be taken as an expression of opinion as to whether or not they are subject to proprietary rights.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Managing Director: Welmoed Spahr

Lead Editor: Celestin Suresh John

Technical Reviewer: Manoj R. Patil

Editorial Board: Steve Anglin, Pramila Balan, Laura Berendson, Aaron Black,

Louise Corrigan, Jonathan Gennick, Robert Hutchinson, Celestin Suresh John,

Nikhil Karkal, James Markham, Susan McDermott, Matthew Moodie, Natalie Pao,

Gwenan Spearing

Coordinating Editor: Prachi Mehta

Copy Editor: Kim Wimpsett

Compositor: SPi Global

Indexer: SPi Global

Artist: SPi Global

Distributed to the book trade worldwide by Springer Science+Business Media New York, 233 Spring Street, 6th Floor, New York, NY 10013. Phone 1-800-SPRINGER, fax (201) 348-4505, e-mail [orders-ny@springer-sbm.com](mailto:orders-ny@springer-sbm.com), or visit [www.springeronline.com](http://www.springeronline.com). Apress Media, LLC is a California LLC and the sole member (owner) is Springer Science + Business Media Finance Inc (SSBM Finance Inc). SSBM Finance Inc is a **Delaware** corporation.

For information on translations, please e-mail [rights@apress.com](mailto:rights@apress.com), or visit [www.apress.com](http://www.apress.com).

Apress and friends of ED books may be purchased in bulk for academic, corporate, or promotional use. eBook versions and licenses are also available for most titles. For more information, reference our Special Bulk Sales–eBook Licensing web page at [www.apress.com/bulk-sales](http://www.apress.com/bulk-sales).

Any source code or other supplementary materials referenced by the author in this text are available to readers at [www.apress.com](http://www.apress.com). For detailed information about how to locate your book's source code, go to [www.apress.com/source-code/](http://www.apress.com/source-code/). Readers can also access source code at SpringerLink in the Supplementary Material section for each chapter.

Printed on acid-free paper

*The six most important people in my life:*  
*The late Kammari Rangaswamy (Teacher)*  
*The late Niranjanamma (Mother)*  
*Devaiah (Father)*  
*Radha (Wife)*  
*Sai Nirupam (Son)*  
*Nitya Maithreyi (Daughter)*

# Contents at a Glance

<b>About the Author .....</b>	<b>xix</b>
<b>About the Technical Reviewer .....</b>	<b>xxi</b>
<b>Acknowledgments .....</b>	<b>xxiii</b>
<b>■ Chapter 1: MapReduce and Its Abstractions .....</b>	<b>1</b>
<b>■ Chapter 2: Data Types.....</b>	<b>21</b>
<b>■ Chapter 3: Grunt .....</b>	<b>33</b>
<b>■ Chapter 4: Pig Latin Fundamentals .....</b>	<b>41</b>
<b>■ Chapter 5: Joins and Functions.....</b>	<b>69</b>
<b>■ Chapter 6: Creating and Scheduling Workflows Using Apache Oozie.....</b>	<b>89</b>
<b>■ Chapter 7: HCatalog.....</b>	<b>103</b>
<b>■ Chapter 8: Pig Latin in Hue.....</b>	<b>115</b>
<b>■ Chapter 9: Pig Latin Scripts in Apache Falcon .....</b>	<b>123</b>
<b>■ Chapter 10: Macros .....</b>	<b>137</b>
<b>■ Chapter 11: User-Defined Functions.....</b>	<b>147</b>
<b>■ Chapter 12: Writing Eval Functions .....</b>	<b>157</b>
<b>■ Chapter 13: Writing Load and Store Functions.....</b>	<b>171</b>
<b>■ Chapter 14: Troubleshooting .....</b>	<b>187</b>
<b>■ Chapter 15: Data Formats.....</b>	<b>201</b>

■ CONTENTS AT A GLANCE

- **Chapter 16: Optimization..... 209**
- **Chapter 17: Hadoop Ecosystem Tools..... 225**
- **Appendix A: Built-in Functions..... 249**
- **Appendix B: Apache Pig in Apache Ambari..... 257**
- **Appendix C: HBaseStorage and ORCStorage Options ..... 261**
- Index..... 265**

# Contents

<b>About the Author .....</b>	<b>xix</b>
<b>About the Technical Reviewer .....</b>	<b>xxi</b>
<b>Acknowledgments .....</b>	<b>xxiii</b>
<b>■ Chapter 1: MapReduce and Its Abstractions .....</b>	<b>1</b>
Small Data Processing .....	1
Relational Database Management Systems .....	3
Data Warehouse Systems .....	3
Parallel Computing .....	4
GFS and MapReduce .....	4
Apache Hadoop .....	4
Problems with MapReduce .....	13
Cascading .....	13
Apache Hive.....	15
Apache Pig.....	16
Summary .....	20
<b>■ Chapter 2: Data Types.....</b>	<b>21</b>
Simple Data Types .....	22
int .....	22
long.....	22
float .....	22
double .....	23
chararray .....	23

boolean .....	23
bytearray .....	23
datetime .....	23
biginteger .....	24
bigdecimal .....	24
Summary of Simple Data Types .....	24
<b>Complex Data Types .....</b>	<b>24</b>
map .....	25
tuple .....	26
bag .....	26
Summary of Complex Data Types .....	27
<b>Schema .....</b>	<b>28</b>
<b>Casting .....</b>	<b>28</b>
Casting Error .....	29
<b>Comparison Operators .....</b>	<b>29</b>
<b>Identifiers .....</b>	<b>30</b>
<b>Boolean Operators .....</b>	<b>31</b>
<b>Summary .....</b>	<b>31</b>
<b>■ Chapter 3: Grunt .....</b>	<b>33</b>
<b>Invoking the Grunt Shell .....</b>	<b>33</b>
<b>Commands .....</b>	<b>34</b>
The fs Command .....	34
The sh Command .....	35
<b>Utility Commands .....</b>	<b>36</b>
help .....	36
history .....	36
quit .....	36
kill .....	37



set.....	37
clear.....	38
exec.....	38
run.....	39
Summary of Commands.....	39
Auto-completion.....	40
Summary.....	40
<b>■ Chapter 4: Pig Latin Fundamentals .....</b>	<b>41</b>
Running Pig Latin Code .....	41
Grunt Shell.....	41
Pig -e .....	42
Pig -f .....	42
Embed Pig Code in a Java Program.....	42
Hue .....	44
Pig Operators and Commands.....	44
Load.....	45
store .....	47
dump .....	48
version.....	48
Foreach Generate .....	48
filter .....	50
Limit.....	51
Assert .....	51
SPLIT.....	52
SAMPLE .....	53
FLATTEN.....	53
import.....	54
define.....	54
distinct.....	55

■ CONTENTS

RANK.....	55
Union .....	56
ORDER BY .....	57
GROUP .....	59
Stream .....	61
MAPREDUCE .....	62
CUBE.....	63
Parameter Substitution .....	65
-param.....	65
-paramfile.....	66
Summary.....	67
<b>■ Chapter 5: Joins and Functions.....</b>	<b>69</b>
Join Operators.....	70
Equi Joins .....	70
cogroup .....	72
CROSS .....	73
Functions.....	74
String Functions .....	74
Mathematical Functions .....	76
Date Functions.....	78
EVAL Functions .....	80
Complex Data Type Functions.....	81
Load/Store Functions.....	82
Summary.....	87
<b>■ Chapter 6: Creating and Scheduling Workflows Using Apache Oozie.....</b>	<b>89</b>
Types of Oozie Jobs.....	89
Workflow.....	89

Using a Pig Latin Script as Part of a Workflow .....	91
Writing job.properties .....	91
workflow.xml .....	91
Uploading Files to HDFS .....	93
Submit the Oozie Workflow .....	93
Scheduling a Pig Script .....	94
Writing the job.properties File .....	94
Writing coordinator.xml .....	94
Upload Files to HDFS .....	96
Submitting Coordinator.....	96
Bundle .....	96
oozie pig Command.....	96
Command-Line Interface.....	98
Job Submitting, Running, and Suspending.....	98
Killing Job.....	98
Retrieving Logs .....	98
Information About a Job .....	98
Oozie User Interface .....	99
Developing Oozie Applications Using Hue .....	100
Summary .....	100
<b>■ Chapter 7: HCatalog.....</b>	<b>103</b>
Features of HCatalog.....	103
Command-Line Interface.....	104
show Command.....	105
Data Definition Language Commands .....	105
dfs and set Commands.....	106

- WebHCatalog ..... 107**
  - Executing Pig Latin Code ..... 108
  - Running a Pig Latin Script from a File ..... 108
  - HCatLoader Example ..... 109
  - Writing the Job Status to a Directory ..... 109
- HCatLoader and HCatStorer ..... 110**
  - Reading Data from HCatalog ..... 110
  - Writing Data to HCatalog ..... 110
  - Running Code ..... 111
  - Data Type Mapping ..... 112
- Summary ..... 113**
- Chapter 8: Pig Latin in Hue ..... 115**
  - Pig Module ..... 115**
    - My Scripts..... 116
    - Pig Helper ..... 117
    - Auto-suggestion ..... 117
    - UDF Usage in Script..... 118
    - Query History..... 118
  - File Browser ..... 119**
  - Job Browser ..... 121**
  - Summary ..... 122**
- Chapter 9: Pig Latin Scripts in Apache Falcon ..... 123**
  - cluster ..... 124**
    - Interfaces..... 124
    - Locations ..... 125
  - feed ..... 126**
    - Feed Types..... 126
    - Frequency..... 126

Late Arrival.....	127
Cluster .....	127
<b>process .....</b>	<b>128</b>
cluster.....	128
Failures.....	128
feed.....	129
workflow.....	129
<b>CLI .....</b>	<b>129</b>
entity .....	129
<b>Web Interface .....</b>	<b>130</b>
Search .....	131
Create an Entity .....	131
Notifications .....	131
Mirror.....	131
<b>Data Replication Using the Falcon Web UI.....</b>	<b>131</b>
Create Cluster Entities .....	132
Create Mirror Job .....	132
<b>Pig Scripts in Apache Falcon.....</b>	<b>134</b>
Oozie Workflow .....	134
Pig Script .....	135
<b>Summary.....</b>	<b>136</b>
<b>■ Chapter 10: Macros .....</b>	<b>137</b>
Structure .....	137
Macro Use Case .....	138
Macro Types .....	138
Internal Macro .....	139
External Macro .....	140

- dryrun..... 141
- Macro Chaining ..... 141
- Macro Rules ..... 142
  - Define Before Usage..... 142
  - Valid Macro Chaining..... 143
  - No Macro Within Nested Block ..... 143
  - No Grunt Shell Commands..... 143
  - Invisible Relations..... 143
- Macro Examples..... 144
  - Macro Without Input Parameters Is Possible..... 144
  - Macro Without Returning Anything Is Possible..... 144
- Summary ..... 145
- **Chapter 11: User-Defined Functions ..... 147**
  - User-Defined Functions..... 148
    - Java ..... 148
    - JavaScript..... 150
    - Other Languages ..... 152
  - Other Libraries..... 154
    - PiggyBank..... 154
    - Apache DataFu ..... 155
  - Summary ..... 155
- **Chapter 12: Writing Eval Functions ..... 157**
  - MapReduce and Pig Features ..... 157
    - Accessing the Distributed Cache..... 157
    - Accessing Counters..... 158
    - Reporting Progress..... 159
    - Output Schema and Input Schema in UDF..... 159
    - Examples of Output and Input Schemas..... 161

Other EVAL Functions .....	162
Algebraic.....	162
Accumulator .....	168
Filter Functions.....	168
Summary.....	169
<b>■ Chapter 13: Writing Load and Store Functions.....</b>	<b>171</b>
Writing a Load Function .....	171
Loading Metadata.....	174
Improving Loader Performance .....	176
Converting from bytearray.....	176
Pushing Down the Predicate .....	177
Writing a Store Function.....	178
Writing Metadata.....	182
Distributed Cache .....	183
Handling Bad Records .....	184
Accessing the Configuration .....	185
Monitoring the UDF Runtime .....	185
Summary.....	186
<b>■ Chapter 14: Troubleshooting .....</b>	<b>187</b>
Illustrate .....	187
describe.....	188
Dump.....	188
Explain.....	188
Plan Types.....	189
Modes.....	193
Unit Testing.....	195
Error Types .....	197

- Counters ..... 198
- Summary ..... 199
- **Chapter 15: Data Formats..... 201**
  - Compression ..... 201
  - Sequence File ..... 202
  - Parquet..... 203
    - Parquet File Processing Using Apache Pig ..... 204
  - ORC..... 205
    - Index ..... 207
    - ACID ..... 207
    - Predicate Pushdown..... 207
    - Data Types ..... 207
    - Benefits ..... 208
  - Summary ..... 208
- **Chapter 16: Optimization..... 209**
  - Advanced Joins ..... 209
    - Small Files ..... 209
    - User-Defined Join Using the Distributed Cache..... 210
    - Big Keys..... 212
    - Sorted Data..... 212
  - Best Practices ..... 213
    - Choose Your Required Fields Early ..... 213
    - Define the Appropriate Schema..... 213
    - Filter Data ..... 214
    - Store Reusable Data ..... 214
    - Use the Algebraic Interface ..... 214
    - Use the Accumulator Interface ..... 215
    - Compress Intermediate Data ..... 215



Combine Small Inputs.....	215
Prefer a Two-Way Join over Multiway Joins.....	216
<b>Better Execution Engine .....</b>	<b>216</b>
<b>Parallelism.....</b>	<b>216</b>
<b>Job Statistics.....</b>	<b>217</b>
<b>Rules .....</b>	<b>218</b>
Partition Filter Optimizer.....	218
Merge foreach .....	218
Constant Calculator .....	219
<b>Cluster Optimization.....</b>	<b>219</b>
Disk Space .....	219
Separate Setup for Zookeeper.....	220
Scheduler .....	220
Name Node Heap Size .....	220
Other Memory Settings .....	221
<b>Summary.....</b>	<b>222</b>
<b>■ Chapter 17: Hadoop Ecosystem Tools.....</b>	<b>225</b>
<b>Apache Zookeeper.....</b>	<b>225</b>
Terminology .....	225
Applications .....	226
Command-Line Interface .....	227
Four-Letter Commands.....	229
Measuring Time .....	230
<b>Cascading.....</b>	<b>230</b>
Defining a Source .....	230
Defining a Sink .....	232
Pipes.....	233
Types of Operations .....	233

- Apache Spark ..... 237
  - Core ..... 238
  - SQL ..... 240
- Apache Tez ..... 245
- Presto ..... 245
  - Architecture ..... 246
  - Connectors ..... 247
  - Pushdown Operations..... 247
- Summary ..... 247
- **Appendix A: Built-in Functions ..... 249**
- **Appendix B: Apache Pig in Apache Ambari ..... 257**
  - Modifying Properties ..... 258
  - Service Check ..... 258
  - Installing Pig..... 259
    - Pig Status ..... 259
    - Check All Available Services..... 259
  - Summary ..... 260
- **Appendix C: HBaseStorage and ORCStorage Options ..... 261**
  - HBaseStorage..... 261
    - Row-Based Conditions ..... 261
    - Timestamp-Based Conditions..... 262
    - Other Conditions ..... 262
  - OrcStorage ..... 263
- Index ..... 265**

# About the Author



**Balaswamy Vaddeman** is a thinker, blogger, and serious and self-motivated big data evangelist with 10 years of experience in IT and 5 years of experience in the big data space. His big data experience covers multiple areas such as analytical applications, product development, consulting, training, book reviews, hackathons, and mentoring. He has proven himself while delivering analytical applications in the retail, banking, and finance domains in three aspects (development, administration, and architecture) of Hadoop-related technologies. At a startup company, he developed a Hadoop-based product that was used for delivering analytical applications without writing code.

In 2013 Balaswamy won the Hadoop Hackathon event for Hyderabad conducted by Cloudwick Technologies. Being the top contributor at [Stackoverflow.com](https://stackoverflow.com), he helped countless people on big data topics at multiple web sites such as [Stackoverflow.com](https://stackoverflow.com) and [Quora.com](https://quora.com). With so much passion on big data, he became an independent trainer and consultant so he could train hundreds of people and set up big data teams in several companies.

# About the Technical Reviewer



**Manoj R. Patil** is a big data architect at TatvaSoft, an IT services and consulting firm. He has a bachelor's of engineering degree from COEP in Pune, India. He is a proven and highly skilled business intelligence professional with 17 years of information technology experience. He is a seasoned BI and big data consultant with exposure to all the leading platforms such as Java EE, .NET, LAMP, and so on. In addition to authoring a book on Pentaho and big data, he believes in knowledge sharing, keeps himself busy in corporate training, and is a passionate teacher. He can be reached at on Twitter @manojrpatil and at <https://in.linkedin.com/in/manojrpatil> on LinkedIn.

Manoj would like to thank his family, especially his two beautiful daughters, Ayushee and Ananyaa, for their patience during the review process.

# Acknowledgments

Writing a book requires a great team. Fortunately, I had a great team for my first project. I am deeply indebted to them for making this project reality.

I would like to thank the publisher, Apress, for providing this opportunity.

Special thanks to Celestin Suresh John for building confidence in me in the initial stages of this project.

Special thanks to Subha Srikant for your valuable feedback. This project would have not been in this shape without you. In fact, I have learned many things from you that could be useful for my future projects also.

Thank you, Manoj R. Patil, for providing valuable technical feedback. Your contribution added a lot of value to this project.

Thank you, Dinesh Kumar, for your valuable time.

Last but not least, thank you, Prachi Mehta, for your prompt coordination.

# CHAPTER 1



# MapReduce and Its Abstractions

In this chapter, you will learn about the technologies that existed before Apache Hadoop, about how Hadoop has addressed the limitations of those technologies, and about the new developments since Hadoop was released.

Data consists of facts collected for analysis. Every business collects data to understand their business and to take action accordingly. In fact, businesses will fall behind their competition if they do not act upon data in a timely manner. Because the number of applications, devices, and users is increasing, data is growing exponentially. Terabytes and petabytes of data have become the norm. Therefore, you need better data management tools for this large amount of data.

Data can be classified into these three types:

- *Small data*: Data is considered small data if it can be measured in gigabytes.
- *Big data*: Big data is characterized by volume, velocity, and variety. *Volume* refers to the size of data, such as terabytes and more. *Velocity* refers to the age of data, such as real-time, near-real-time, and streaming data. *Variety* talks about types of data; there are mainly three types of data: structured, semistructured, and unstructured.
- *Fast data*: Fast data is a type of big data that is useful for the real-time presentation of data. Because of the huge demand for real-time or near-real-time data, fast data is evolving in a separate and unique space.

## Small Data Processing

Many tools and technologies are available for processing small data. You can use languages such as Python, Perl, and Java, and you can use relational database management systems (RDBMSs) such as Oracle, MySQL, and Postgres. You can even use data warehousing tools and extract/transform/load (ETL) tools. In this section, I will discuss how small data processing is done.

---

**Electronic supplementary material** The online version of this chapter (doi:[10.1007/978-1-4842-2337-6\\_1](https://doi.org/10.1007/978-1-4842-2337-6_1)) contains supplementary material, which is available to authorized users.

Assume you have the following text in a file called *fruits*:

```
Apple, grape
Apple, grape, pear
Apple, orange
```

Let's write a program in a shell script that first filters out the word *pear* and then counts the number of words in the file. Here's the code:

```
cat fruits|tr ',' '\n'|grep -v -i 'pear'|sort -f|uniq -c -i
```

This code is explained in the following paragraphs.

In this code, *tr* (for “translate” or “transliterate”) is a Unix program that takes two inputs and replaces the first set of characters with the second set of characters. In the previous program, the *tr* program replaces each comma (,) with a new line character (\n). *grep* is a command used for searching for specific text. So, the previous program performs an inverse search on the word *pear* using the *-v* option and ignores the case using *-i*.

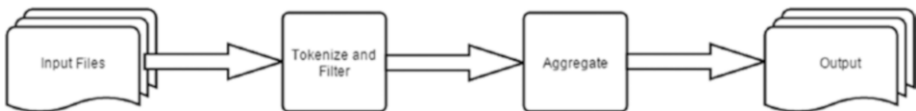
The *sort* command produces data in sorted order. The *-f* option ignores case while sorting.

*uniq* is a Unix program that combines adjacent lines from the input file for reporting purposes. In the previous program, *uniq* takes sorted words from the *sort* command output and generates the word count. The *-c* option is for the count, and the *-i* option is for ignoring case.

The program produces the following output:

```
Apple 3
Grape 2
Orange 1
```

You can divide program functionality into two stages; first is tokenize and filtering, and second is aggregation. *Sort* is supporting functionality of aggregation. Figure 1-1 shows the program flow.



**Figure 1-1.** Program flow

The previous program can be run on a single machine and on small data. Such simple programs can be used to perform simple operations such as searching and sorting on one file at a time. However, writing complex queries involving multiple files and multiple conditions requires better data processing tools. Database management systems (DBMS) and RDBMS technologies were developed to address querying problems with structured data.

## Relational Database Management Systems

RDBMSs were developed based on the relational model founded by E. F. Codd. There are many commercial RDBMS products such as Oracle, SQL Server, and DB2. Many open source RDBMSs such as MySQL, Postgres, and SQLite are also popular. RDBMSs store data in tables, and you can define relations between tables.

Here are some advantages of RDBMSs:

- RDBMS products come with sophisticated query languages that can easily retrieve data from multiple tables with multiple conditions.
- The query language used in RDBMSs is called Structured Query Language (SQL); it provides easy data definition, manipulation, and control.
- RDBMSs also support transactions.
- RDBMSs support low-latency queries so users can access databases interactively, and they are also useful for online transaction processing (OLTP).

RDBMSs have these disadvantages:

- As data is stored in table format, RDBMSs support only structured data.
- You need to define a schema at the time of loading data.
- RDBMSs can scale only to gigabytes of data, and they are mainly designed for frequent updates.

Because the data size in today's organizations has grown exponentially, RDBMSs have not been able to scale with respect to data size. Processing terabytes of data can take days.

Having terabytes of data has become the norm for almost all businesses. And new data types like semistructured and unstructured have arrived. Semistructured data has a partial structure like in web server log files, and it needs to be parsed like Extensible Markup Language (XML) in order to analyze it. Unstructured data does not have any structure; this includes images, videos, and e-books.

## Data Warehouse Systems

Data warehouse systems were introduced to address the problems of RDBMSs. Data warehouse systems such as Teradata are able to scale up to terabytes of data, and they are mainly used for OLAP use cases.

Data warehousing systems have these disadvantages:

- Data warehouse systems are a costly solution.
- They still cannot process other data types such as semistructured and unstructured data.
- They cannot scale to petabytes and beyond.



All traditional data-processing technologies experience a couple of common problems: storage and performance.

Computing infrastructure can face the problem of node failures. Data needs to be available irrespective of node failures, and storage systems should be able to store large volumes of data.

Traditional data processing technologies used a scale-up approach to process a large volume of data. A scale-up approach adds more computing power to existing nodes, so it cannot scale to petabytes and more because the rest of computing infrastructure becomes a performance bottleneck.

Growing storage and processing needs have created a need for new technologies such as parallel computing technologies.

## Parallel Computing

The following are several parallel computing technologies.

### GFS and MapReduce

Google has created two parallel computing technologies to address the storage and processing problems of big data. They are Google File System (GFS) and MapReduce. GFS is a distributed file system that provides fault tolerance and high performance on commodity hardware. GFS follows a master-slave architecture. The master is called Master, and the slave is called ChunkServer in GFS. MapReduce is an algorithm based on key-value pairs used for processing a huge amount of data on commodity hardware. These are two successful parallel computing technologies that address the storage and processing limitations of big data.

### Apache Hadoop

Apache Hadoop is an open source framework used for storing and processing large data sets on commodity hardware in a fault-tolerant manner.

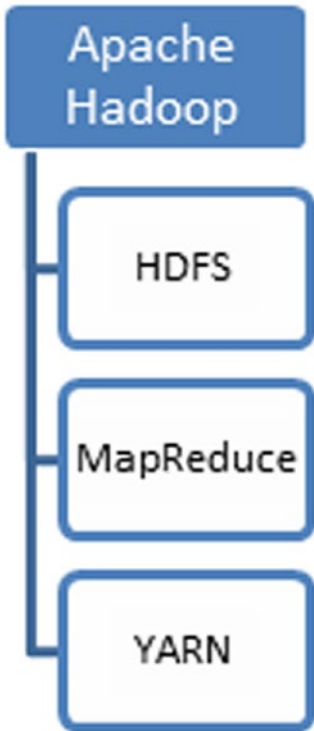
Hadoop was written by Doug Cutting and Mark Cafarella in 2006 while working for Yahoo to improve the performance of the Nutch search engine. Cutting named it after his son's stuffed elephant toy. In 2007, it was given to the Apache Software Foundation.

Initially Hadoop was adopted by Yahoo and, later, by companies like Facebook and Microsoft. Yahoo has about 100,000 CPUs and 40,000 nodes for Hadoop. The largest Hadoop cluster has about 4,500 nodes. Yahoo runs about **850,000 Hadoop jobs** every day. Unlike conventional parallel computing technologies, Hadoop follows a scale-out strategy, which makes it more scalable. In fact, Apache Hadoop had set a benchmark by sorting 1.42 **terabytes per minute**.

Most of Hadoop is written in Java, but it has support for many programming languages such as C, C++, Python, and Scala through its streaming module. Apache Hadoop was initially written for high throughput and batch-processing systems. RDBMS technologies were written for frequent modifications in data, whereas Hadoop has been written for frequent reads.

Moore's law says the processing capability of a machine will double every two years. Kryder's law says the storage capacity of disks will grow faster than Moore's law. The cost of computing and storage devices will go down every day, and these two factors can support more scalable technologies. Apache Hadoop was designed while keeping these things in mind, and parallel computing technologies like this will become more common going forward.

The latest Apache Hadoop contains three modules, as shown in Figure 1-2. They are HDFS, MapReduce, and Yet Another Resource Negotiator (YARN).



**Figure 1-2.** *The three components of Hadoop*

## HDFS

The Hadoop distributed file system is used for storing large data sets. It divides files into blocks and stores every block on at least multiple nodes. This is called a *replication factor*, and by default it is 3. HDFS is fault-tolerant because it has more than one replica for every block, so it can handle node failures without affecting data processing. A block of HDFS is the same as an operating system block, but a HDFS block size is larger, such as 64 MB or 128 MB. Unlike traditional storage systems, it is highly scalable. It does not require any special hardware and can work on commodity hardware.

Assume you have a replication factor of 3, a block size of 64 MB, and 640 MB of data needs to be uploaded into HDFS. At the time of uploading the data into HDFS, 640 MB is divided into 10 blocks with respect to block size. Every block is stored on three nodes, which would consume 1920 MB of space on a cluster.

HDFS follows a master-slave architecture. The master is called the *name node*, and the slave is called a *data node*. The data node is fault tolerant because the same block is replicated to two more nodes. The name node was a single point of failure in initial versions; in fact, Hadoop used to go down if the name node crashed. But Hadoop 2.0+ versions have high availability of the name node. If the active name node is down, the standby name node becomes active without affecting the running jobs.

## MapReduce

MapReduce is key-value programming model used for processing large data sets. It has two core functions: Map and Reduce. They are derived from functional programming languages. Both functions take a key-value pair as input and generate a key-value pair as output.

The Map task is responsible for filtering operations and preparing the data required for the Reduce tasks. The Map task will generate intermediate output and write it to the hard disk. For every key that is being generated by the Map task, a Reduce node is identified and will be sent to the key for further processing.

The Map task takes the key-value pair as input and generates the key-value pair as output.

*(key1, value1) -----> Map Task-----> (Key2, Value2)*

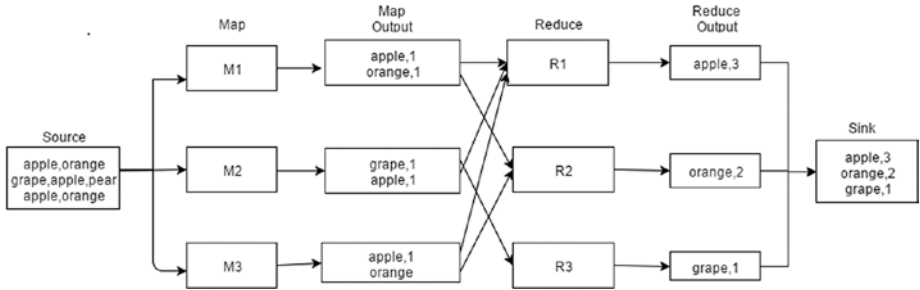
The Reduce task is responsible for data aggregation operations such as count, max, min, average, and so on. A reduce operation will be performed on a per-key basis. Every functionality can be expressed in MapReduce.

The Reduce task takes the key and list of values as input and generates the key and value as output.

*(key2, List (value2))-----> Reduce Task -----> (Key3, value3)*

In addition to the Map and Reduce tasks, there is an extra stage called the *combiner* to improve the performance of MapReduce. The combiner will do partial aggregation on the Map side so that the Map stage has to write less data to disk.

You will now see how MapReduce generates a word count. Figure 1-3 depicts how MapReduce generates the fruits word count after filtering out the word *pear*.



**Figure 1-3.** MapReduce generating a word count

Source and Sink are HDFS directories. When you upload data to HDFS, data is divided into chunks called *blocks*. Blocks will be processed in a parallel manner on all available nodes.

The first stage is Map, which performs filtering and data preparation after tokenization. All Map tasks (M1, M2, and M3) will do the initial numbering for words that are useful for the final aggregation. And M2 filters out the word *pear*.

The key and list of its values are retrieved from the Map output and sent to the reducer node. For example, the Apple key and its values (1, 1, 1) are sent to the reducer node R1. The reducer aggregates all values to generate the count output.

Between Map and Reduce, there is an internal stage called *shuffling* where the reducer node for the map output is identified.

You will now see how to write the same word count program using MapReduce. You first need to write a mapper class for the Map stage.

## Writing a Map Class

The following is the Map program that is used for the same tokenization and data filtering as in the shell script discussed earlier:

```
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;
public class FilterMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
    private final static IntWritable one = new IntWritable(1);
    private Text word = new Text();
    public void map(LongWritable offset, Text line, Context context) throws
        IOException, InterruptedException {
        //tokenize line with comma as delimiter
        StringTokenizer itr = new StringTokenizer(line.toString(),",");
        //Iterate all tokens and filter pear word
        while (itr.hasMoreTokens()) {
```

```

        String strToken=itr.nextToken();
        if(!strToken.equals("pear")){
//converting string data type to text data type of mapreduce

        word.set(strToken);
        context.write(word, one);//Map output
        }
    }
}

```

The Map class should extend the Mapper class, which has parameters for the input key, input value, output key, and output value. You need to override the map() method. This code specifies LongWritable for the input key, Text for the input value, Text for the output key, and IntWritable for the output value.

In the map() method, you use StringTokenizer to convert a sentence into words. You are iterating words using a while loop, and you are filtering the word pear using an if loop. The Map stage output is written to context.

For every run of the map() method, the line offset value is the input key, the line is the input value, the word in the line will become an output key, and 1 is the output value, as shown in Figure 1-4.



**Figure 1-4.** M2 stage

The map() method runs once per every line. It tokenizes the line into words, and it filters the word pear before writing other words with the default of 1.

If the combiner is available, the combiner is run before the Reduce stage. Every Map task will have a combiner task that will produce aggregated output. Assume you have two apple words in the second line that is processed by the M2 map task.

The Map output without the combiner will look like Figure 1-5.



**Figure 1-5.** Map output without the combiner