

---

# Professional Hadoop Solutions

---

Hadoop Real-World Solutions Cookbook Second Edition  
Hadoop Application Architectures  
Securing Hadoop  
Big Data Analytics with Hadoop 3  
Mastering Hadoop 3  
Modern Big Data Processing with Hadoop  
Big Data Made Easy  
Handbook of Research on Big Data Storage and Visualization Techniques  
Professional Hadoop  
Scaling Big Data with Hadoop and Solr  
Pro Microsoft HDInsight  
Apache Hadoop 3 Quick Start Guide  
Practical Hive  
Hadoop Real World Solutions Cookbook  
Data Algorithms  
Pro Hadoop Data Analytics  
Optimizing Hadoop for MapReduce  
Apache Hive Essentials  
Professional Hadoop  
Professional Hadoop Solutions  
Understanding Big Data: Analytics for Enterprise Class Hadoop and Streaming Data  
Field Guide to Hadoop  
Hadoop Security  
Hadoop Essentials  
BIG DATA ANALYTICS AND HADOOP SOLUTIONS WITH SAS TOOLS  
Big Data Networked Storage Solution for Hadoop  
Big Data Forensics - Learning Hadoop Investigations  
Beginning Apache Hadoop Administration  
Data Analytics with Hadoop  
Practical Hadoop Migration  
Professional Hadoop  
Architecting Big Data: Mastering Hadoop Solution  
Hadoop: Data Processing and Modelling  
Big Data Processing With Hadoop  
Apache Hadoop YARN  
Hadoop: The Definitive Guide  
Pro Hadoop  
Hadoop Real-World Solutions Cookbook  
Hadoop For Dummies  
Architecting Modern Data Platforms

---

**LIA BURGESS**


---

**Hadoop Real-World Solutions Cookbook Second Edition** Packt Publishing Ltd

Ready to use statistical and machine-learning techniques across large data sets? This practical guide shows you why the Hadoop ecosystem is perfect for the job. Instead of deployment, operations, or software development usually associated with distributed computing, you'll focus on particular analyses you can build, the data warehousing techniques that Hadoop provides, and higher order data workflows this framework can produce. Data scientists and analysts will learn how to perform a wide range of techniques, from writing MapReduce and Spark applications with Python to using advanced modeling and data management with Spark MLlib, Hive, and HBase. You'll also learn about the analytical processes and data systems available to build and empower data products that can handle—and actually require—huge amounts of data. Understand core concepts behind Hadoop and cluster computing Use design patterns and

parallel analytical algorithms to create distributed data analysis jobs Learn about data management, mining, and warehousing in a distributed context using Apache Hive and HBase Use Sqoop and Apache Flume to ingest data from relational databases Program complex Hadoop and Spark applications with Apache Pig and Spark DataFrames Perform machine learning techniques such as classification, clustering, and collaborative filtering with Spark's MLlib [Hadoop Application Architectures](#) "O'Reilly Media, Inc."

Big Data represents a new era in data exploration and utilization, and IBM is uniquely positioned to help clients navigate this transformation. This book reveals how IBM is leveraging open source Big Data technology, infused with IBM technologies, to deliver a robust, secure, highly available, enterprise-class Big Data platform. The three defining characteristics of Big Data—volume, variety, and velocity—are discussed. You'll get a primer on Hadoop and how IBM is hardening it for the enterprise, and learn when to leverage

IBM InfoSphere BigInsights (Big Data at rest) and IBM InfoSphere Streams (Big Data in motion) technologies. Industry use cases are also included in this practical guide. Learn how IBM hardens Hadoop for enterprise-class scalability and reliability Gain insight into IBM's unique in-motion and at-rest Big Data analytics platform Learn tips and tricks for Big Data use cases and solutions Get a quick Hadoop primer [Securing Hadoop](#) "O'Reilly Media, Inc."

A comprehensive guide to mastering the most advanced Hadoop 3 concepts Key Features Get to grips with the newly introduced features and capabilities of Hadoop 3 Crunch and process data using MapReduce, YARN, and a host of tools within the Hadoop ecosystem Sharpen your Hadoop skills with real-world case studies and code Book Description Apache Hadoop is one of the most popular big data solutions for distributed storage and for processing large chunks of data. With Hadoop 3, Apache promises to provide a high-performance, more fault-tolerant, and highly efficient big data

processing platform, with a focus on improved scalability and increased efficiency. With this guide, you'll understand advanced concepts of the Hadoop ecosystem tool. You'll learn how Hadoop works internally, study advanced concepts of different ecosystem tools, discover solutions to real-world use cases, and understand how to secure your cluster. It will then walk you through HDFS, YARN, MapReduce, and Hadoop 3 concepts. You'll be able to address common challenges like using Kafka efficiently, designing low latency, reliable message delivery Kafka systems, and handling high data volumes. As you advance, you'll discover how to address major challenges when building an enterprise-grade messaging system, and how to use different stream processing systems along with Kafka to fulfil your enterprise goals. By the end of this book, you'll have a complete understanding of how components in the Hadoop ecosystem are effectively integrated to implement a fast and reliable data pipeline, and you'll be equipped to tackle a range of real-world problems in data

pipelines. What you will learnGain an in-depth understanding of distributed computing using Hadoop 3Develop enterprise-grade applications using Apache Spark, Flink, and moreBuild scalable and high-performance Hadoop data pipelines with security, monitoring, and data governanceExplore batch data processing patterns and how to model data in HadoopMaster best practices for enterprises using, or planning to use, Hadoop 3 as a data platformUnderstand security aspects of Hadoop, including authorization and authenticationWho this book is for If you want to become a big data professional by mastering the advanced concepts of Hadoop, this book is for you. You'll also find this book useful if you're a Hadoop professional looking to strengthen your knowledge of the Hadoop ecosystem. Fundamental knowledge of the Java programming language and basics of Hadoop is necessary to get started with this book. *Big Data Analytics with Hadoop 3* IGI Global Due to the increasing availability of affordable internet services, the

number of users, and the need for a wider range of multimedia-based applications, internet usage is on the rise. With so many users and such a large amount of data, the requirements of analyzing large data sets leads to the need for further advancements to information processing. Big Data Processing With Hadoop is an essential reference source that discusses possible solutions for millions of users working with a variety of data applications, who expect fast turnaround responses, but encounter issues with processing data at the rate it comes in. Featuring research on topics such as market basket analytics, scheduler load simulator, and writing YARN applications, this book is ideally designed for IoT professionals, students, and engineers seeking coverage on many of the real-world challenges regarding big data. *Mastering Hadoop 3* "O'Reilly Media, Inc." The go-to guidebook for deploying Big Data solutions with Hadoop Today's enterprise architects need to understand how the Hadoop frameworks and APIs fit together, and how

they can be integrated to deliver real-world solutions. This book is a practical, detailed guide to building and implementing those solutions, with code-level instruction in the popular Wrox tradition. It covers storing data with HDFS and Hbase, processing data with MapReduce, and automating data processing with Oozie. Hadoop security, running Hadoop with Amazon Web Services, best practices, and automating Hadoop processes in real time are also covered in depth. With in-depth code examples in Java and XML and the latest on recent additions to the Hadoop ecosystem, this complete resource also covers the use of APIs, exposing their inner workings and allowing architects and developers to better leverage and customize them. The ultimate guide for developers, designers, and architects who need to build and deploy Hadoop applications

Covers storing and processing data with various technologies, automating data processing, Hadoop security, and delivering real-time solutions

Includes detailed, real-world examples and code-level guidelines Explains

when, why, and how to use these tools effectively

Written by a team of Hadoop experts in the programmer-to-programmer Wrox style

Professional Hadoop Solutions is the reference enterprise architects and developers need to maximize the power of Hadoop.

Modern Big Data Processing with Hadoop

Packt Publishing Ltd

You've heard the hype about Hadoop: it runs petabyte-scale data mining tasks insanely fast, it runs gigantic tasks on clouds for absurdly cheap, it's been heavily committed to by tech giants like IBM, Yahoo!, and the Apache Project, and it's completely open-source (thus free). But what exactly is it, and more importantly, how do you even get a Hadoop cluster up and running?

From Apress, the name you've come to trust for hands-on technical knowledge, Pro Hadoop brings you up to speed on Hadoop. You learn the ins and outs of MapReduce; how to structure a cluster, design, and implement the Hadoop file system; and how to build your first cloud-computing tasks using Hadoop. Learn how to let Hadoop take care of distributing and

parallelizing your software—you just focus on the code, Hadoop takes care of the rest. Best of all, you'll learn from a tech professional who's been in the Hadoop scene since day one. Written from the perspective of a principal engineer with down-in-the-trenches knowledge of what to do wrong with Hadoop, you learn how to avoid the common, expensive first errors that everyone makes with creating their own Hadoop system or inheriting someone else's. Skip the novice stage and the expensive, hard-to-fix mistakes...go straight to seasoned pro on the hottest cloud-computing framework with Pro Hadoop. Your productivity will blow your managers away.

*Big Data Made Easy*

Anand Vemula

Call Data Record Analytics using Hive

**Handbook of Research on Big Data Storage and Visualization**

**Techniques** Apress

Ready to unlock the power of your data? With this comprehensive guide, you'll learn how to build and maintain reliable, scalable, distributed systems with Apache Hadoop. This book is ideal for programmers looking

to analyze datasets of any size, and for administrators who want to set up and run Hadoop clusters. You'll find illuminating case studies that demonstrate how Hadoop is used to solve specific problems. This third edition covers recent changes to Hadoop, including material on the new MapReduce API, as well as MapReduce 2 and its more flexible execution model (YARN). Store large datasets with the Hadoop Distributed File System (HDFS) Run distributed computations with MapReduce Use Hadoop's data and I/O building blocks for compression, data integrity, serialization (including Avro), and persistence Discover common pitfalls and advanced features for writing real-world MapReduce programs Design, build, and administer a dedicated Hadoop cluster—or run Hadoop in the cloud Load data from relational databases into HDFS, using Sqoop Perform large-scale data processing with the Pig query language Analyze datasets with Hive, Hadoop's data warehousing system Take advantage of HBase for structured and semi-structured data, and

ZooKeeper for building distributed systems  
**Professional Hadoop** John Wiley & Sons  
 This book is a step-by-step tutorial filled with practical examples which will focus mainly on the key security tools and implementation techniques of Hadoop security. This book is great for Hadoop practitioners (solution architects, Hadoop administrators, developers, and Hadoop project managers) who are looking to get a good grounding in what Kerberos is all about and who wish to learn how to implement end-to-end Hadoop security within an enterprise setup. It's assumed that you will have some basic understanding of Hadoop as well as be familiar with some basic security concepts.  
*Scaling Big Data with Hadoop and Solr* John Wiley & Sons  
 Realistic, simple code examples to solve problems at scale with Hadoop and related technologies  
**Pro Microsoft HDInsight** Apress  
 Explore big data concepts, platforms, analytics, and their applications using the power of Hadoop 3  
 Key Features Learn Hadoop 3 to build

effective big data analytics solutions on-premise and on cloud Integrate Hadoop with other big data tools such as R, Python, Apache Spark, and Apache Flink Exploit big data using Hadoop 3 with real-world examples Book Description Apache Hadoop is the most popular platform for big data processing, and can be combined with a host of other big data tools to build powerful analytics solutions. Big Data Analytics with Hadoop 3 shows you how to do just that, by providing insights into the software as well as its benefits with the help of practical examples. Once you have taken a tour of Hadoop 3's latest features, you will get an overview of HDFS, MapReduce, and YARN, and how they enable faster, more efficient big data processing. You will then move on to learning how to integrate Hadoop with the open source tools, such as Python and R, to analyze and visualize data and perform statistical computing on big data. As you get acquainted with all this, you will explore how to use Hadoop 3 with Apache Spark and Apache Flink for real-time data analytics and stream

processing. In addition to this, you will understand how to use Hadoop to build analytics solutions on the cloud and an end-to-end pipeline to perform big data analysis using practical use cases. By the end of this book, you will be well-versed with the analytical capabilities of the Hadoop ecosystem. You will be able to build powerful solutions to perform big data analytics and get insight effortlessly. What you will learn Explore the new features of Hadoop 3 along with HDFS, YARN, and MapReduce Get well-versed with the analytical capabilities of Hadoop ecosystem using practical examples Integrate Hadoop with R and Python for more efficient big data processing Learn to use Hadoop with Apache Spark and Apache Flink for real-time data analytics Set up a Hadoop cluster on AWS cloud Perform big data analytics on AWS using Elastic Map Reduce Who this book is for Big Data Analytics with Hadoop 3 is for you if you are looking to build high-performance analytics solutions for your enterprise or business using Hadoop 3's powerful features, or you're new to big data analytics. A basic understanding of the Java

programming language is required.

*Apache Hadoop 3 Quick Start Guide* Packt Publishing Ltd

This book is an example-based tutorial that deals with Optimizing Hadoop for MapReduce job performance. If you are a Hadoop administrator, developer, MapReduce user, or beginner, this book is the best choice available if you wish to optimize your clusters and applications. Having prior knowledge of creating MapReduce applications is not necessary, but will help you better understand the concepts and snippets of MapReduce class template code.

*Practical Hive* "O'Reilly Media, Inc."

A comprehensive guide to design, build and execute effective Big Data strategies using Hadoop Key Features -Get an in-depth view of the Apache Hadoop ecosystem and an overview of the architectural patterns pertaining to the popular Big Data platform - Conquer different data processing and analytics challenges using a multitude of tools such as Apache Spark, Elasticsearch, Tableau and more -A comprehensive, step-by-

step guide that will teach you everything you need to know, to be an expert Hadoop Architect Book Description The complex structure of data these days requires sophisticated solutions for data transformation, to make the information more accessible to the users. This book empowers you to build such solutions with relative ease with the help of Apache Hadoop, along with a host of other Big Data tools. This book will give you a complete understanding of the data lifecycle management with Hadoop, followed by modeling of structured and unstructured data in Hadoop. It will also show you how to design real-time streaming pipelines by leveraging tools such as Apache Spark, and build efficient enterprise search solutions using Elasticsearch. You will learn to build enterprise-grade analytics solutions on Hadoop, and how to visualize your data using tools such as Apache Superset. This book also covers techniques for deploying your Big Data solutions on the cloud Apache Ambari, as well as expert techniques for managing and administering your Hadoop cluster. By the

end of this book, you will have all the knowledge you need to build expert Big Data systems. What you will learn Build an efficient enterprise Big Data strategy centered around Apache Hadoop Gain a thorough understanding of using Hadoop with various Big Data frameworks such as Apache Spark, Elasticsearch and more Set up and deploy your Big Data environment on premises or on the cloud with Apache Ambari Design effective streaming data pipelines and build your own enterprise search solutions Utilize the historical data to build your analytics solutions and visualize them using popular tools such as Apache Superset Plan, set up and administer your Hadoop cluster efficiently Who this book is for This book is for Big Data professionals who want to fast-track their career in the Hadoop industry and become an expert Big Data architect. Project managers and mainframe professionals looking forward to build a career in Big Data Hadoop will also find this book to be useful. Some understanding of Hadoop is required to get the best out of this book.

Hadoop Real World Solutions Cookbook Packt Publishing Ltd  
Learn advanced analytical techniques and leverage existing tool kits to make your analytic applications more powerful, precise, and efficient. This book provides the right combination of architecture, design, and implementation information to create analytical systems that go beyond the basics of classification, clustering, and recommendation. Pro Hadoop Data Analytics emphasizes best practices to ensure coherent, efficient development. A complete example system will be developed using standard third-party components that consist of the tool kits, libraries, visualization and reporting code, as well as support glue to provide a working and extensible end-to-end system. The book also highlights the importance of end-to-end, flexible, configurable, high-performance data pipeline systems with analytical components as well as appropriate visualization results. You'll discover the importance of mix-and-match or hybrid systems, using different analytical components in one application. This hybrid

approach will be prominent in the examples. What You'll Learn Build big data analytic systems with the Hadoop ecosystem Use libraries, tool kits, and algorithms to make development easier and more effective Apply metrics to measure performance and efficiency of components and systems Connect to standard relational databases, noSQL data sources, and more Follow case studies with example components to create your own systems Who This Book Is For Software engineers, architects, and data scientists with an interest in the design and implementation of big data analytical systems using Hadoop, the Hadoop ecosystem, and other associated technologies. **Data Algorithms** Notion Press  
Let Hadoop For Dummies help harness the power of your data and rein in the information overload Big data has become big business, and companies and organizations of all sizes are struggling to find ways to retrieve valuable information from their massive data sets with becoming overwhelmed. Enter Hadoop and this easy-to-understand For Dummies guide. Hadoop

For Dummies helps readers understand the value of big data, make a business case for using Hadoop, navigate the Hadoop ecosystem, and build and manage Hadoop applications and clusters. Explains the origins of Hadoop, its economic benefits, and its functionality and practical applications Helps you find your way around the Hadoop ecosystem, program MapReduce, utilize design patterns, and get your Hadoop cluster up and running quickly and easily Details how to use Hadoop applications for data mining, web analytics and personalization, large-scale text processing, data science, and problem-solving Shows you how to improve the value of your Hadoop cluster, maximize your investment in Hadoop, and avoid common pitfalls when building your Hadoop cluster From programmers challenged with building and maintaining affordable, scaleable data systems to administrators who must deal with huge volumes of information effectively and efficiently, this how-to has something to help you with Hadoop.

[Pro Hadoop Data Analytics](#)  
Apress

Many corporations are finding that the size of their data sets are outgrowing the capability of their systems to store and process them. The data is becoming too big to manage and use with traditional tools. The solution: implementing a big data system. As *Big Data Made Easy: A Working Guide to the Complete Hadoop Toolset* shows, Apache Hadoop offers a scalable, fault-tolerant system for storing and processing data in parallel. It has a very rich toolset that allows for storage (Hadoop), configuration (YARN and ZooKeeper), collection (Nutch and Solr), processing (Storm, Pig, and Map Reduce), scheduling (Oozie), moving (Sqoop and Avro), monitoring (Chukwa, Ambari, and Hue), testing (Big Top), and analysis (Hive). The problem is that the Internet offers IT pros wading into big data many versions of the truth and some outright falsehoods born of ignorance. What is needed is a book just like this one: a wide-ranging but easily understood set of instructions to explain where to get Hadoop tools, what they can do, how to install them, how to configure them, how to

integrate them, and how to use them successfully. And you need an expert who has worked in this area for a decade—someone just like author and big data expert Mike Frampton. *Big Data Made Easy* approaches the problem of managing massive data sets from a systems perspective, and it explains the roles for each project (like architect and tester, for example) and shows how the Hadoop toolset can be used at each system stage. It explains, in an easily understood manner and through numerous examples, how to use each tool. The book also explains the sliding scale of tools available depending upon data size and when and how to use them. *Big Data Made Easy* shows developers and architects, as well as testers and project managers, how to: Store big data Configure big data Process big data Schedule processes Move data among SQL and NoSQL systems Monitor data Perform big data analytics Report on big data processes and projects Test big data systems *Big Data Made Easy* also explains the best part, which is that this toolset is free.

Anyone can download it and—with the help of this book—start to use it within a day. With the skills this book will teach you under your belt, you will add value to your company or client immediately, not to mention your career.

[Optimizing Hadoop for MapReduce](#) Packt Publishing Ltd

This IBM® Redpaper™ provides a reference architecture, based on Apache Hadoop, to help businesses gain control over their data, meet tight service level agreements (SLAs) around their data applications, and turn data-driven insight into effective action. Big Data Networked Storage Solution for Hadoop delivers the capabilities for ingesting, storing, and managing large data sets with high reliability. IBM InfoSphere® Big Insights™ provides an innovative analytics platform that processes and analyzes all types of data to turn large complex data into insight. IBM InfoSphere BigInsights brings the power of Hadoop to the enterprise. With built-in analytics, extensive integration capabilities, and the reliability, security and support that you require, IBM can help

put your big data to work for you. This IBM Redpaper publication provides basic guidelines and best practices for how to size and configure Big Data Networked Storage Solution for Hadoop.

**Apache Hive Essentials**  
Apress

Dive into the world of SQL on Hadoop and get the most out of your Hive data warehouses. This book is your go-to resource for using Hive: authors Scott Shaw, Ankur Gupta, David Kjerrumgaard, and Andreas Francois Vermeulen take you through learning HiveQL, the SQL-like language specific to Hive, to analyze, export, and massage the data stored across your Hadoop environment. From deploying Hive on your hardware or virtual machine and setting up its initial configuration to learning how Hive interacts with Hadoop, MapReduce, Tez and other big data technologies, Practical Hive gives you a detailed treatment of the software. In addition, this book discusses the value of open source software, Hive performance tuning, and how to leverage semi-structured and unstructured data. What

You Will Learn Install and configure Hive for new and existing datasets  
Perform DDL operations  
Execute efficient DML operations  
Use tables, partitions, buckets, and user-defined functions  
Discover performance tuning tips and Hive best practices  
Who This Book Is For Developers, companies, and professionals who deal with large amounts of data and could use software that can efficiently manage large volumes of input. It is assumed that readers have the ability to work with SQL.

**Professional Hadoop**  
Apress

Over 90 hands-on recipes to help you learn and master the intricacies of Apache Hadoop 2.X, YARN, Hive, Pig, Oozie, Flume, Sqoop, Apache Spark, and Mahout  
About This Book- Implement outstanding Machine Learning use cases on your own analytics models and processes.- Solutions to common problems when working with the Hadoop ecosystem.- Step-by-step implementation of end-to-end big data use cases.  
Who This Book Is For Readers who have a basic knowledge of big data systems and want to

advance their knowledge with hands-on recipes. **What You Will Learn**- Installing and maintaining Hadoop 2.X cluster and its ecosystem.- Write advanced Map Reduce programs and understand design patterns.- **Advanced Data Analysis** using the Hive, Pig, and Map Reduce programs.- Import and export data from various sources using Sqoop and Flume.- Data storage in various file formats such as Text, Sequential, Parquet, ORC, and RC Files.- Machine learning principles with libraries such as Mahout- Batch and Stream data processing using Apache Spark  
**In Detail** Big data is the current requirement. Most organizations produce huge amount of data every day. With the arrival of Hadoop-like tools, it has become easier for everyone to solve big data problems with great efficiency and at minimal cost. Grasping Machine Learning techniques will help you greatly in building predictive models and using this data to make the right decisions for your organization. **Hadoop Real World Solutions Cookbook** gives readers insights into learning and mastering big data via

recipes. The book not only clarifies most big data tools in the market but also provides best practices for using them. The book provides recipes that are based on the latest versions of Apache Hadoop 2.X, YARN, Hive, Pig, Sqoop, Flume, Apache Spark, Mahout and many more such ecosystem tools. This real-world-solution cookbook is packed with handy recipes you can apply to your own everyday issues. Each chapter provides in-depth recipes that can be referenced easily. This book provides detailed practices on the latest technologies such as YARN and Apache Spark. Readers will be able to consider themselves as big data experts on completion of this book. This guide is an invaluable tutorial if you are planning to implement a big data warehouse for your business. **Style and approach** An easy-to-follow guide that walks you through world of big data. Each tool in the Hadoop ecosystem is explained in detail and the recipes are placed in such a manner that readers can implement them sequentially. Plenty of reference links are provided for advanced reading.

**Professional Hadoop Solutions** Createspace Independent Publishing Platform

If you are ready to dive into the MapReduce framework for processing large datasets, this practical book takes you step by step through the algorithms and tools you need to build distributed MapReduce applications with Apache Hadoop or Apache Spark. Each chapter provides a recipe for solving a massive computational problem, such as building a recommendation system. You'll learn how to implement the appropriate MapReduce solution with code that you can use in your projects. Dr. Mahmoud Parsian covers basic design patterns, optimization techniques, and data mining and machine learning solutions for problems in bioinformatics, genomics, statistics, and social network analysis. This book also includes an overview of MapReduce, Hadoop, and Spark. **Topics include:** Market basket analysis for a large set of transactions Data mining algorithms (K-means, KNN, and Naive Bayes) Using huge genomic data to sequence DNA and RNA Naive Bayes

theorem and Markov  
chains for data and  
market prediction  
Recommendation  
algorithms and pairwise

document similarity  
Linear regression, Cox  
regression, and Pearson  
correlation Allelic  
frequency and mining

DNA Social network  
analysis (recommendation  
systems, counting  
triangles, sentiment  
analysis)