
Apache Hadoop 3 0 0 Hdfs Architecture

Thank you very much for downloading **Apache Hadoop 3 0 0 Hdfs Architecture**. Maybe you have knowledge that, people have see numerous time for their favorite books gone this Apache Hadoop 3 0 0 Hdfs Architecture, but end taking place in harmful downloads.

Rather than enjoying a fine book past a cup of coffee in the afternoon, instead they juggled taking into consideration some harmful virus inside their computer. **Apache Hadoop 3 0 0 Hdfs Architecture** is within reach in our digital library an online entrance to it is set as public as a result you can download it instantly. Our digital library saves in multiple countries, allowing you to acquire the most less latency period to download any of our books considering this one. Merely said, the Apache Hadoop 3 0 0 Hdfs Architecture is universally compatible following any devices to read.

*Apache
Hadoop 3 0 0
Hdfs
Architecture* *Downloaded
from
ftp.wagmtv.com
by guest*

WALSH RANDOLPH

Parallel Computing

*Technologies Packt
Publishing Ltd
Pro Microsoft HDInsight
is a complete guide to
deploying and using
Apache Hadoop on the*

Microsoft Windows Azure Platforms. The information in this book enables you to process enormous volumes of structured as well as non-structured data easily using HDInsight, which is Microsoft's own distribution of Apache Hadoop. Furthermore, the blend of Infrastructure as a Service (IaaS) and Platform as a Service (PaaS) offerings available through Windows Azure lets you take advantage of Hadoop's processing power without the worry of creating, configuring, maintaining, or managing your own cluster. With the data explosion that is soon to happen, the open source Apache Hadoop Framework is gaining traction, and it benefits

from a huge ecosystem that has risen around the core functionalities of the Hadoop distributed file system (HDFS™) and Hadoop Map Reduce. Pro Microsoft HDInsight equips you with the knowledge, confidence, and technique to configure and manage this ecosystem on Windows Azure. The book is an excellent choice for anyone aspiring to be a data scientist or data engineer, putting you a step ahead in the data mining field. Guides you through installation and configuration of an HDInsight cluster on Windows Azure Provides clear examples of configuring and executing Map Reduce jobs Helps you consume data and

diagnose errors from the Windows Azure HDInsight Service

A Working Guide to the Complete Hadoop Toolset 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.

Elasticsearch 5.x

Cookbook Packt Publishing Ltd
Web Scraping techniques are getting more popular, since data is as valuable as oil in 21st century. Through this book get some key knowledge about using XPath, regEX; web scraping libraries for R like rvest and RSelenium technologies. Key Features Techniques, tools and frameworks for web scraping with R Scrape data effortlessly from a variety of websites Learn how to selectively choose the data to scrape, and build your dataset
Book Description Web scraping is a technique

to extract data from websites. It simulates the behavior of a website user to turn the website itself into a web service to retrieve or introduce new data. This book gives you all you need to get started with scraping web pages using R programming. You will learn about the rules of RegEx and Xpath, key components for scraping website data. We will show you web scraping techniques, methodologies, and frameworks. With this book's guidance, you will become comfortable with the tools to write and test RegEx and XPath rules. We will focus on examples of dynamic websites for scraping data and how to implement the techniques learned. You will learn how to

collect URLs and then create XPath rules for your first web scraping script using rvest library. From the data you collect, you will be able to calculate the statistics and create R plots to visualize them. Finally, you will discover how to use Selenium drivers with R for more sophisticated scraping. You will create AWS instances and use R to connect a PostgreSQL database hosted on AWS. By the end of the book, you will be sufficiently confident to create end-to-end web scraping systems using R. What you will learn

Write and create regEX rules Write XPath rules to query your data Learn how web scraping methods work Use rvest to crawl web pages Store data retrieved from the web

Learn the key uses of R Selenium to scrape data Who this book is for This book is for R programmers who want to get started quickly with web scraping, as well as data analysts who want to learn scraping using R. Basic knowledge of R is all you need to get started with this book.

A Classroom Approach

Apress

Over 150 recipes to design and optimize large scale Apache Cassandra deployments.

Hadoop on Windows

Springer Nature

This book constitutes the thoroughly refereed post-conference

proceedings of the 8th

TPC Technology

Conference, on

Performance

Evaluation and

Benchmarking, TPCTC 2016, held in conjunction with the 41st International Conference on Very Large Databases (VLDB 2016) in New Delhi, India, in September 2016. The 9 papers presented were carefully reviewed and selected from 20 submissions. They reflect the rapid pace at which industry experts and researchers develop innovative techniques for evaluation, measurement and characterization of complex systems.

Hands-On for Developers and Technical Professionals
"O'Reilly Media, Inc."
Advanced analytics on your Big Data with latest Apache Spark 2.x About This Book An advanced guide with a combination of

instructions and practical examples to extend the most up-to-date Spark functionalities. Extend your data processing capabilities to process huge chunk of data in minimum time using advanced concepts in Spark. Master the art of real-time processing with the help of Apache Spark 2.x Who This Book Is For If you are a developer with some experience with Spark and want to strengthen your knowledge of how to get around in the world of Spark, then this book is ideal for you. Basic knowledge of Linux, Hadoop and Spark is assumed. Reasonable knowledge of Scala is expected. What You Will Learn Examine Advanced Machine Learning and DeepLearning with

MLlib, SparkML, SystemML, H2O and DeepLearning4J Study highly optimised unified batch and real-time data processing using SparkSQL and Structured Streaming Evaluate large-scale Graph Processing and Analysis using GraphX and GraphFrames Apply Apache Spark in Elastic deployments using Jupyter and Zeppelin Notebooks, Docker, Kubernetes and the IBM Cloud Understand internal details of cost based optimizers used in Catalyst, SystemML and GraphFrames Learn how specific parameter settings affect overall performance of an Apache Spark cluster Leverage Scala, R and python for your data science projects In Detail Apache Spark is

an in-memory cluster-based parallel processing system that provides a wide range of functionalities such as graph processing, machine learning, stream processing, and SQL. This book aims to take your knowledge of Spark to the next level by teaching you how to expand Spark's functionality and implement your data flows and machine/deep learning programs on top of the platform. The book commences with an overview of the Spark ecosystem. It will introduce you to Project Tungsten and Catalyst, two of the major advancements of Apache Spark 2.x. You will understand how memory management and binary processing, cache-aware computation, and code

generation are used to speed things up dramatically. The book extends to show how to incorporate H2O, SystemML, and DeepLearning4j for machine learning, and Jupyter Notebooks and Kubernetes/Docker for cloud-based Spark. During the course of the book, you will learn about the latest enhancements to Apache Spark 2.x, such as interactive querying of live data and unifying DataFrames and Datasets. You will also learn about the updates on the APIs and how DataFrames and Datasets affect SQL, machine learning, graph processing, and streaming. You will learn to use Spark as a big data operating system, understand how to implement advanced analytics on

the new APIs, and explore how easy it is to use Spark in day-to-day tasks. Style and approach This book is an extensive guide to Apache Spark modules and tools and shows how Spark's functionality can be extended for real-time processing and storage with worked examples. [Techniques and tools to crawl and scrape data from websites](#) Packt Publishing Ltd Bigdata is one of the most demanding markets in the IT sector. If you are an administrator or a have a passion for knowing the internal configurations of Hadoop, then this book is for you. This book enables a professional to learn about Hadoop in terms of installation, configuration, and management. This

book will help the reader to jumpstart with Hadoop frameworks, its ecosystem components and slowly progress towards learning the administration part of Hadoop. The level of this book goes from beginner to intermediate with 70% hands-on exercises. Some of the techniques that you will learn include, • Installation and configuration of Hadoop cluster • Performing Hadoop Cluster Upgrade • Understanding and implementing HDFS Federation • Understanding and Implementing High Availability • Implementing HA on a Federated Cluster • Zookeeper CLI • Apache Hive Installation and

Security • HBase Multi-master setup • Oozie installation, configuration and job submission • Setting up HDFS Quotas • Setting up HDFS NFS gateway • Understanding and implementing rolling upgrade and much more.

Real-Time Analytics
Pragmatic Bookshelf
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.

Recipes for Scaling Up with Hadoop and Spark
Packt Publishing Ltd
This book constitutes the thoroughly refereed post-conference proceedings of the 8th TPC Technology Conference, on Performance Evaluation and Benchmarking, TPCTC 2017, held in conjunction with the 43rd International Conference on Very Large Databases (VLDB 2017) in August/September 2017. The 12 papers presented were carefully reviewed and selected from numerous submissions. The TPC remains committed to developing new benchmark standards to keep pace with these rapid changes in technology.

Beginning Apache Hadoop Administration Guru99

The authors provide an understanding of big data and MapReduce by clearly presenting the basic terminologies and concepts. They have employed over 100 illustrations and many worked-out examples to convey the concepts and methods used in big data, the inner workings of MapReduce, and single node/multi-node installation on physical/virtual machines. This book covers almost all the necessary information on Hadoop MapReduce for most online certification exams. Upon completing this book, readers will find it easy to understand other big data processing tools such

as Spark, Storm, etc. Ultimately, readers will be able to:

- understand what big data is and the factors that are involved
- understand the inner workings of MapReduce, which is essential for certification exams
- learn the features and weaknesses of MapReduce
- set up Hadoop clusters with 100s of physical/virtual machines
- create a virtual machine in AWS
- write MapReduce with Eclipse in a simple way
- understand other big data processing tools and their applications

Pro MongoDB Development Springer
Build efficient data flow and machine learning programs with this flexible, multi-functional open-source cluster-computing

framework Key Features Master the art of real-time big data processing and machine learning Explore a wide range of use-cases to analyze large data Discover ways to optimize your work by using many features of Spark 2.x and Scala Book Description Apache Spark is an in-memory, cluster-based data processing system that provides a wide range of functionalities such as big data processing, analytics, machine learning, and more. With this Learning Path, you can take your knowledge of Apache Spark to the next level by learning how to expand Spark's functionality and building your own data flow and machine learning programs on this platform. You will

work with the different modules in Apache Spark, such as interactive querying with Spark SQL, using DataFrames and datasets, implementing streaming analytics with Spark Streaming, and applying machine learning and deep learning techniques on Spark using MLlib and various external tools. By the end of this elaborately designed Learning Path, you will have all the knowledge you need to master Apache Spark, and build your own big data processing and analytics pipeline quickly and without any hassle. This Learning Path includes content from the following Packt products: Mastering Apache Spark 2.x by Romeo Kienzler Scala and Spark for Big Data

Analytics by Md. Rezaul Karim, Sridhar Alla Apache Spark 2.x Machine Learning Cookbook by Siamak Amirghodsi, Meenakshi Rajendran, Broderick Hall, Shuen Mei Cookbook What you will learn Get to grips with all the features of Apache Spark 2.x Perform highly optimized real-time big data processing Use ML and DL techniques with Spark MLlib and third-party tools Analyze structured and unstructured data using SparkSQL and GraphX Understand tuning, debugging, and monitoring of big data applications Build scalable and fault-tolerant streaming applications Develop scalable recommendation engines Who this book is for If you are an

intermediate-level Spark developer looking to master the advanced capabilities and use-cases of Apache Spark 2.x, this Learning Path is ideal for you. Big data professionals who want to learn how to integrate and use the features of Apache Spark and build a strong big data pipeline will also find this Learning Path useful. To grasp the concepts explained in this Learning Path, you must know the fundamentals of Apache Spark and Scala. Techniques to Analyze and Visualize Streaming Data Springer Nature A fast paced guide that will help you learn about Apache Hadoop 3 and its ecosystem Key Features Set up,

configure and get started with Hadoop to get useful insights from large data sets. Work with the different components of Hadoop such as MapReduce, HDFS and YARN. Learn about the new features introduced in Hadoop 3.

Book Description

Apache Hadoop is a widely used distributed data platform. It enables large datasets to be efficiently processed instead of using one large computer to store and process the data. This book will get you started with the Hadoop ecosystem, and introduce you to the main technical topics, including MapReduce, YARN, and HDFS. The book begins with an overview of big data and Apache Hadoop. Then, you will set up a pseudo

Hadoop development environment and a multi-node enterprise Hadoop cluster. You will see how the parallel programming paradigm, such as MapReduce, can solve many complex data processing problems. The book also covers the important aspects of the big data software development lifecycle, including quality assurance and control, performance, administration, and monitoring. You will then learn about the Hadoop ecosystem, and tools such as Kafka, Sqoop, Flume, Pig, Hive, and HBase. Finally, you will look at advanced topics, including real time streaming using Apache Storm, and data analytics using Apache Spark. By the end of the book, you

will be well versed with different configurations of the Hadoop 3 cluster. What you will learn Store and analyze data at scale using HDFS, MapReduce and YARN Install and configure Hadoop 3 in different modes Use Yarn effectively to run different applications on Hadoop based platform Understand and monitor how Hadoop cluster is managed Consume streaming data using Storm, and then analyze it using Spark Explore Apache Hadoop ecosystem components, such as Flume, Sqoop, HBase, Hive, and Kafka Who this book is for Aspiring Big Data professionals who want to learn the essentials of Hadoop 3 will find this book to be useful. Existing Hadoop

users who want to get up to speed with the new features introduced in Hadoop 3 will also benefit from this book. Having knowledge of Java programming will be an added advantage. *Advanced Analytics with Spark* Springer Hadoop has changed the way large data sets are analyzed, stored, transferred, and processed. At such low cost, it provides benefits like supports partial failure, fault tolerance, consistency, scalability, flexible schema, and so on. It also supports cloud computing. More and more number of individuals are looking forward to mastering their Hadoop skills. While initiating with Hadoop, most users are unsure about how to proceed with

Hadoop. They are not aware of what are the pre-requisite or data structure they should be familiar with. Or How to make the most efficient use of Hadoop and its ecosystem. To help them with all these queries and other issues this e-book is designed. The book gives insights into many of Hadoop libraries and packages that are not known to many Big data Analysts and Architects. The e-book also tells you about Hadoop MapReduce and HDFS. The example in the e-book is well chosen and demonstrates how to control Hadoop ecosystem through various shell commands. With this book, users will gain expertise in Hadoop technology and its related components.

The book leverages you with the best Hadoop content with the lowest price range. After going through this book, you will also acquire knowledge on Hadoop Security required for Hadoop Certifications like CCAH and CCDH. It is a definite guide to Hadoop. Table Of Content Chapter 1: What Is Big Data 1. Examples Of 'Big Data' 2. Categories Of 'Big Data' 3. Characteristics Of 'Big Data' 4. Advantages Of Big Data Processing Chapter 2: Introduction to Hadoop 1. Components of Hadoop 2. Features Of 'Hadoop' 3. Network Topology In Hadoop Chapter 3: Hadoop Installation Chapter 4: HDFS 1. Read Operation 2. Write Operation 3. Access HDFS using

JAVA API 4. Access HDFS Using COMMAND-LINE INTERFACE Chapter 5: Mapreduce 1. How MapReduce works 2. How MapReduce Organizes Work? Chapter 6: First Program 1. Understanding MapReducer Code 2. Explanation of SalesMapper Class 3. Explanation of SalesCountryReducer Class 4. Explanation of SalesCountryDriver Class Chapter 7: Counters & Joins In MapReduce 1. Two types of counters 2. MapReduce Join Chapter 8: MapReduce Hadoop Program To Join Data Chapter 9: Flume and Sqoop 1. What is SQOOP in Hadoop? 2. What is FLUME in Hadoop? 3. Some Important features of FLUME Chapter 10: Pig 1. Introduction to PIG 2. Create your First PIG Program 3. PART 1) Pig Installation 4. PART 2) Pig Demo Chapter 11: OOZIE 1. What is OOZIE? 2. How does OOZIE work? 3. Example Workflow Diagram 4. Oozie workflow application 5. Why use Oozie? 6. FEATURES OF OOZIE *Large-Scale Graph Processing Using Apache Giraph* "O'Reilly Media, Inc." Navigate the world of data analysis, visualization, and machine learning with over 100 hands-on Scala recipes About This Book Implement Scala in your data analysis using features from Spark, Breeze, and Zeppelin Scale up your data analytics infrastructure with practical recipes for

Scala machine learning Recipes for every stage of the data analysis process, from reading and collecting data to distributed analytics

Who This Book Is For This book shows data scientists and analysts how to leverage their existing knowledge of Scala for quality and scalable data analysis.

What You Will Learn Familiarize and set up the Breeze and Spark libraries and use data structures Import data from a host of possible sources and create dataframes from CSV Clean, validate and transform data using Scala to pre-process numerical and string data Integrate quintessential machine learning algorithms using Scala stack Bundle and scale up Spark jobs by deploying them into a

variety of cluster managers Run streaming and graph analytics in Spark to visualize data, enabling exploratory analysis In Detail This book will introduce you to the most popular Scala tools, libraries, and frameworks through practical recipes around loading, manipulating, and preparing your data. It will also help you explore and make sense of your data using stunning and insightful visualizations, and machine learning toolkits. Starting with introductory recipes on utilizing the Breeze and Spark libraries, get to grips with how to import data from a host of possible sources and how to pre-process numerical, string, and date data. Next, you'll get an

understanding of concepts that will help you visualize data using the Apache Zeppelin and Bokeh bindings in Scala, enabling exploratory data analysis. iscover how to program quintessential machine learning algorithms using Spark ML library. Work through steps to scale your machine learning models and deploy them into a standalone cluster, EC2, YARN, and Mesos. Finally dip into the powerful options presented by Spark Streaming, and machine learning for streaming data, as well as utilizing Spark GraphX. Style and approach This book contains a rich set of recipes that covers the full spectrum of interesting data analysis tasks and will

help you revolutionize your data analysis skills using Scala and Spark.
Apache Hadoop 3 Quick Start Guide
O'Reilly Media
Scala will be a valuable tool to have on hand during your data science journey for everything from data cleaning to cutting-edge machine learning
About This Book Build data science and data engineering solutions with ease An in-depth look at each stage of the data analysis process — from reading and collecting data to distributed analytics Explore a broad variety of data processing, machine learning, and genetic algorithms through diagrams, mathematical formulations, and source code Who This

Book Is For This learning path is perfect for those who are comfortable with Scala programming and now want to enter the field of data science. Some knowledge of statistics is expected. What You Will Learn Transfer and filter tabular data to extract features for machine learning Read, clean, transform, and write data to both SQL and NoSQL databases Create Scala web applications that couple with JavaScript libraries such as D3 to create compelling interactive visualizations Load data from HDFS and HIVE with ease Run streaming and graph analytics in Spark for exploratory analysis Bundle and scale up Spark jobs by deploying them into a variety of cluster

managers Build dynamic workflows for scientific computing Leverage open source libraries to extract patterns from time series Master probabilistic models for sequential data In Detail Scala is especially good for analyzing large sets of data as the scale of the task doesn't have any significant impact on performance. Scala's powerful functional libraries can interact with databases and build scalable frameworks — resulting in the creation of robust data pipelines. The first module introduces you to Scala libraries to ingest, store, manipulate, process, and visualize data. Using real world examples, you will learn how to design

scalable architecture to process and model data — starting from simple concurrency constructs and progressing to actor systems and Apache Spark. After this, you will also learn how to build interactive visualizations with web frameworks. Once you have become familiar with all the tasks involved in data science, you will explore data analytics with Scala in the second module. You'll see how Scala can be used to make sense of data through easy to follow recipes. You will learn about Bokeh bindings for exploratory data analysis and quintessential machine learning with algorithms with Spark ML library. You'll get a sufficient

understanding of Spark streaming, machine learning for streaming data, and Spark graphX. Armed with a firm understanding of data analysis, you will be ready to explore the most cutting-edge aspect of data science — machine learning. The final module teaches you the A to Z of machine learning with Scala. You'll explore Scala for dependency injections and implicits, which are used to write machine learning algorithms. You'll also explore machine learning topics such as clustering, dimensionality reduction, Naive Bayes, Regression models, SVMs, neural networks, and more. This learning path combines some of the best that Packt has to

offer into one complete, curated package. It includes content from the following Packt products: Scala for Data Science, Pascal Bugnion Scala Data Analysis Cookbook, Arun Manivannan Scala for Machine Learning, Patrick R. Nicolas Style and approach A complete package with all the information necessary to start building useful data engineering and data science solutions straight away. It contains a diverse set of recipes that cover the full spectrum of interesting data analysis tasks and will help you revolutionize your data analysis skills using Scala.

Big Data Made Easy
Apress

This book constitutes the proceedings of the

14th International Conference on Parallel Computing Technologies, PaCT 2017, held in Nizhny Novgorod, Russia, in September 2017. The 25 full papers and 24 short papers presented were carefully reviewed and selected from 93 submissions. The papers are organized in topical sections on mainstream parallel computing, parallel models and algorithms in numerical computation, cellular automata and discrete event systems, organization of parallel computation, parallel computing applications.

15th IFIP TC8 International Conference, CISIM 2016, Vilnius, Lithuania, September 14-16,

2016, Proceedings

Packt Publishing Ltd
Learn how to use the Apache Hadoop projects, including MapReduce, HDFS, Apache Hive, Apache HBase, Apache Kafka, Apache Mahout, and Apache Solr. From setting up the environment to running sample applications each chapter in this book is a practical tutorial on using an Apache Hadoop ecosystem project. While several books on Apache Hadoop are available, most are based on the main projects, MapReduce and HDFS, and none discusses the other Apache Hadoop ecosystem projects and how they all work together as a cohesive big data development platform. What You Will Learn: Set up the

environment in Linux for Hadoop projects using Cloudera Hadoop Distribution CDH 5 Run a MapReduce job Store data with Apache Hive, and Apache HBase Index data in HDFS with Apache Solr Develop a Kafka messaging system Stream Logs to HDFS with Apache Flume Transfer data from MySQL database to Hive, HDFS, and HBase with Sqoop Create a Hive table over Apache Solr Develop a Mahout User Recommender System Who This Book Is For: Apache Hadoop developers. Pre-requisite knowledge of Linux and some knowledge of Hadoop is required.

Moving Beyond MapReduce and Batch Processing with Apache Hadoop 2 Notion Press
This book constitutes

the proceedings of the 15th IFIP TC8 International Conference on Computer Information Systems and Industrial Management, CISIM 2016, held in Vilnius, Lithuania, in September 2016. The 63 regular papers presented together with 1 invited paper and 5 keynotes in this volume were carefully reviewed and selected from about 89 submissions. The main topics covered are rough set methods for big data analytics; images, visualization, classification; optimization, tuning; scheduling in manufacturing and other applications; algorithms; decisions; intelligent distributed systems; and biometrics, identification, security.

Learn about big data processing and analytics Springer

Pro MongoDB Development is about MongoDB, a NoSQL database based on the BSON (binary JSON) document model. The book discusses all aspects of using MongoDB in web applications: Java, PHP, Ruby, JavaScript are the most commonly used programming/scripting languages and the book discusses accessing MongoDB database with these languages. The book also discusses using Java EE frameworks Kundera and Spring Data with MongoDB. As NoSQL databases are commonly used with the Hadoop ecosystem the book also discusses using MongoDB with Apache Hive. Migration

from other NoSQL databases (Apache Cassandra and Couchbase) and from relational databases (Oracle Database) is also discussed. What You'll Learn: How to use a Java client and MongoDB shell How to use MongoDB with PHP, Ruby, and Node.js as well How to migrate Apache Cassandra tables to MongoDB documents; Couchbase to MongoDB; and transferring data between Oracle and MongoDB How to use Kundera, Spring Data, and Spring XD with MongoDB How to load MongoDB data into Oracle Database and integrating MongoDB with Oracle Database in Oracle Data Integrator Audience: The target audience of the book is NoSQL database developers.

Target audience includes Java, PHP and Ruby developers. The book is suitable for an intermediate level course in NoSQL database. Mastering Apache Spark 2.x Packt Publishing Ltd Data is bigger, arrives faster, and comes in a variety of formats—and it all needs to be processed at scale for analytics or machine learning. But how can you process such varied workloads efficiently? Enter Apache Spark. Updated to include Spark 3.0, this second edition shows data engineers and data scientists why structure and unification in Spark matters. Specifically, this book explains how to perform simple and complex data analytics and employ machine

learning algorithms. Through step-by-step walk-throughs, code snippets, and notebooks, you'll be able to: Learn Python, SQL, Scala, or Java high-level Structured APIs Understand Spark operations and SQL Engine Inspect, tune, and debug Spark operations with Spark configurations and Spark UI Connect to

data sources: JSON, Parquet, CSV, Avro, ORC, Hive, S3, or Kafka Perform analytics on batch and streaming data using Structured Streaming Build reliable data pipelines with open source Delta Lake and Spark Develop machine learning pipelines with MLib and productionize models using MLflow