

## Hadoop The Definitive Guide

Hadoop: The Definitive Guide helps you harness the power of your data. Ideal for processing large datasets, the Apache Hadoop framework is an open source implementation of the MapReduce algorithm on which Google built its empire. This comprehensive resource demonstrates how to use Hadoop to build reliable, scalable, distributed systems: programmers will find details for analyzing large datasets, and administrators will learn how to set up and run Hadoop clusters.

Discover how Apache Hadoop can unleash the power of your data. This comprehensive resource shows you how to build and maintain reliable, scalable, distributed systems with the Hadoop framework -- an open source implementation of MapReduce, the algorithm on which Google built its empire. Programmers will find details for analyzing datasets of any size, and administrators will learn how to set up and run Hadoop clusters. This revised edition covers recent changes to Hadoop, including new features such as Hive, Sqoop, and Avro. It also provides illuminating case studies that illustrate how Hadoop is used to solve specific problems. Looking to get the most out of your data? This is your book. Use the Hadoop Distributed File System (HDFS) for storing large datasets, then run distributed computations over those datasets with MapReduce Become familiar with Hadoop's data and I/O building blocks for compression, data integrity, serialization, 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 Use Pig, a high-level query language for large-scale data processing Analyze datasets with Hive, Hadoop's data warehousing system Take advantage of HBase, Hadoop's database for structured and semi-structured data Learn ZooKeeper, a toolkit of coordination primitives for building distributed systems "Now you have the opportunity to learn about Hadoop from a master -- not only of the technology, but also of common sense and plain talk." --Doug Cutting, Cloudera

Written by Ganglia designers and maintainers, this book shows you how to collect and visualize metrics from clusters, grids, and cloud infrastructures at any scale. Want to track CPU utilization from 50,000 hosts every ten seconds? Ganglia is just the tool you need, once you know how its main components work together. This hands-on book helps experienced system administrators take advantage of Ganglia 3.x. Learn how to extend the base set of metrics you collect, fetch current values, see aggregate views of metrics, and observe time-series trends in your data. You'll also examine real-world case studies of Ganglia installs that feature challenging monitoring requirements. Determine whether Ganglia is a good fit for your environment Learn how Ganglia's gmond and gmetad daemons build a metric collection overlay Plan for scalability early in your Ganglia deployment, with valuable tips and advice Take data visualization to a new level with gweb, Ganglia's web frontend Write plugins to extend gmond's metric-collection capability Troubleshoot issues you may

encounter with a Ganglia installation Integrate Ganglia with the sFlow and Nagios monitoring systems Contributors include: Robert Alexander, Jeff Buchbinder, Frederiko Costa, Alex Dean, Dave Josephsen, Peter Phaal, and Daniel Pocock. Case study writers include: John Allspaw, Ramon Bastiaans, Adam Compton, Andrew Dibble, and Jonah Horowitz.

Learn how to use, deploy, and maintain Apache Spark with this comprehensive guide, written by the creators of the open-source cluster-computing framework. With an emphasis on improvements and new features in Spark 2.0, authors Bill Chambers and Matei Zaharia break down Spark topics into distinct sections, each with unique goals. You'll explore the basic operations and common functions of Spark's structured APIs, as well as Structured Streaming, a new high-level API for building end-to-end streaming applications. Developers and system administrators will learn the fundamentals of monitoring, tuning, and debugging Spark, and explore machine learning techniques and scenarios for employing MLlib, Spark's scalable machine-learning library. Get a gentle overview of big data and Spark Learn about DataFrames, SQL, and Datasets—Spark's core APIs—through worked examples Dive into Spark's low-level APIs, RDDs, and execution of SQL and DataFrames Understand how Spark runs on a cluster Debug, monitor, and tune Spark clusters and applications Learn the power of Structured Streaming, Spark's stream-processing engine Learn how you can apply MLlib to a variety of problems, including classification or recommendation

This book gathers selected papers presented at the 3rd Conference on Computing Systems and Applications (CSA'2018), held at the Ecole Militaire Polytechnique, Algiers, Algeria on April 24–25, 2018. The CSA'2018 constitutes a leading forum for exchanging, discussing and leveraging modern computer systems technology in such varied fields as: data science, computer networks and security, information systems and software engineering, and computer vision. The contributions presented here will help promote and advance the adoption of computer science technologies in industrial, entertainment, social, and everyday applications. Though primarily intended for students, researchers, engineers and practitioners working in the field, it will also benefit a wider audience interested in the latest developments in the computer sciences.

If you've been asked to maintain large and complex Hadoop clusters, this book is a must. Demand for operations-specific material has skyrocketed now that Hadoop is becoming the de facto standard for truly large-scale data processing in the data center. Eric Sammer, Principal Solution Architect at Cloudera, shows you the particulars of running Hadoop in production, from planning, installing, and configuring the system to providing ongoing maintenance. Rather than run through all possible scenarios, this pragmatic operations guide calls out what works, as demonstrated in critical deployments. Get a high-level overview of HDFS and MapReduce: why they exist and how they work Plan a Hadoop

deployment, from hardware and OS selection to network requirements Learn setup and configuration details with a list of critical properties Manage resources by sharing a cluster across multiple groups Get a runbook of the most common cluster maintenance tasks Monitor Hadoop clusters—and learn troubleshooting with the help of real-world war stories Use basic tools and techniques to handle backup and catastrophic failure

Hadoop: The Definitive Guide helps you harness the power of your data. Ideal for processing large datasets, the Apache Hadoop framework is an open source implementation of the MapReduce algorithm on which Google built its empire. This comprehensive resource demonstrates how to use Hadoop to build reliable, scalable, distributed systems: programmers will find details for analyzing large datasets, and administrators will learn how to set up and run Hadoop clusters.

Complete with case studies that illustrate how Hadoop solves specific problems, this book helps you: Use the Hadoop Distributed File System (HDFS) for storing large datasets, and run distributed computations over those datasets using MapReduce Become familiar with Hadoop's data and I/O building blocks for compression, data integrity, serialization, 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 Use Pig, a high-level query language for large-scale data processing Take advantage of HBase, Hadoop's database for structured and semi-structured data Learn ZooKeeper, a toolkit of coordination primitives for building distributed systems If you have lots of data -- whether it's gigabytes or petabytes -- Hadoop is the perfect solution. Hadoop: The Definitive Guide is the most thorough book available on the subject. "Now you have the opportunity to learn about Hadoop from a master-not only of the technology, but also of common sense and plain talk." -- Doug Cutting, Hadoop Founder, Yahoo!

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.

Internationale bestseller over de impact van technologie op ons leven: Google Glasses, zelfrijdende auto's, computers die het menselijk brein vervangen... De digitalisering heeft ons leven drastisch veranderd, en we staan nog maar aan het begin van deze revolutie. 'Vanaf nu wordt de verandering pas echt duizelingwekkend', aldus Erik Brynjolfsson en Andrew McAfee, beiden verbonden aan het prestigieuze MIT. 'En het is aanpassen of verliezen.' Miljoenen mensen dreigen hun baan te verliezen, precaire machtsevenwichten verschuiven en de sociale ongelijkheid groeit. Dit tweede tijdperk der machines kan echter ook zorgen voor meer welvaart. Maar dan moeten we nu de juiste keuzes maken.

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.

This book encompasses empirical evidences to understand the application of data analytical techniques in emerging contexts. Varied studies relating to manufacturing and services sectors including healthcare, banking, information technology, power, education sector etc. stresses upon the systematic approach followed in applying the data analytical techniques; and also analyses how these techniques are effective in decision-making in different contexts. Especially, the application of regression modeling, financial modelling, multi-group modeling, cluster analysis, and sentiment analysis will help the readers in understanding critical business scenarios in the best possible way, and which later can help them in arriving at best solution for the business related problems. The individual chapters will help the readers in understanding the role of specific data analytic tools and techniques in resolving business operational issues experienced in manufacturing and service organisations in India and in developing countries. The book offers a relevant resource that

will help readers in the application and interpretation of data analytical statistical practices relating to emerging issues like customer experience, marketing capability, quality of manufactured products, strategic orientation, high-performance human resource policy, employee resilience, financial resources, etc. This book will be of interest to a professional audience that include practitioners, policy makers, NGOs, managers and employees as well as academicians, researchers and students.

If you're looking for a scalable storage solution to accommodate a virtually endless amount of data, this updated edition shows you how Apache HBase can meet your needs. Modeled after Google's BigTable architecture, HBase scales to billions of rows and millions of columns, while ensuring that write and read performance remain constant. Fully revised for HBase 1.0, this second edition brings you up to speed on the new HBase client API, as well as security features and new case studies that demonstrate HBase use in the real world. Whether you just started to evaluate this non-relational database, or plan to put it into practice right away, this book has your back. Launch into basic, advanced, and administrative features of HBase's new client-facing API Use new classes to integrate HBase with Hadoop's MapReduce framework Explore HBase's architecture, including the storage format, write-ahead log, and background processes Dive into advanced usage, such extended client and server options Learn cluster sizing, tuning, and monitoring best practices Design schemas, copy tables, import bulk data, decommission nodes, and other tasks Go deeper into HBase security, including Kerberos and encryption at rest

This book presents a focus on proteins and their structures. The text describes various scalable solutions for protein structure similarity searching, carried out at main representation levels and for prediction of 3D structures of proteins. Emphasis is placed on techniques that can be used to accelerate similarity searches and protein structure modeling processes. The content of the book is divided into four parts. The first part provides background information on proteins and their representation levels, including a formal model of a 3D protein structure used in computational processes, and a brief overview of the technologies used in the solutions presented in the book. The second part of the book discusses Cloud services that are utilized in the development of scalable and reliable cloud applications for 3D protein structure similarity searching and protein structure prediction. The third part of the book shows the utilization of scalable Big Data computational frameworks, like Hadoop and Spark, in massive 3D protein structure alignments and identification of intrinsically disordered regions in protein structures. The fourth part of the book focuses on finding 3D protein structure similarities, accelerated with the use of GPUs and the use of multithreading and relational databases for efficient approximate searching on protein secondary structures. The book introduces advanced techniques and computational architectures that benefit from recent achievements in the field of computing and parallelism. Recent developments in

computer science have allowed algorithms previously considered too time-consuming to now be efficiently used for applications in bioinformatics and the life sciences. Given its depth of coverage, the book will be of interest to researchers and software developers working in the fields of structural bioinformatics and biomedical databases.

This book constitutes the thoroughly refereed post-conference proceedings of the First International Workshop on Algorithmic Aspects of Cloud Computing, ALGO CLOUD 2015, held in Patras, Greece, in September 2015 in conjunction with ALGO 2015. The 13 revised full papers presented together with 2 tutorial papers were carefully reviewed and selected from 37 initial submissions. They cover a wide range of topics in two main tracks: algorithmic aspects of large-scale data stores, and software tools and distributed architectures for cloud-based data management.

Onze gegevens worden gebruikt om ons te bespioneren en om ons dingen te verkopen die we niet willen en ook niet nodig hebben. Maar met de enorme hoeveelheid gegevens die we op internet achterlaten ('big data') is nog iets veel interessanter te doen. Ons gedrag online, wanneer we ons onbespied wanen, onthult wie we echt zijn. Als beheerder van een datingsite beschikt Rudder over een schat aan informatie over wat we leuk vinden, met wie we praten, wat we daarbij drinken en hoe laat we naar bed gaan. Het is een nieuwe manier om psychologisch onderzoek te doen, veel effectiever dan de traditionele vragenlijst, waarbij we ons altijd beter voordoen dan we zijn. Het is misschien even slikken, maar Christian Rudder laat zien dat Facebook, Google en OkCupid ons beter kennen dan onze beste vrienden. Rudder is een geestige reisgids door de jungle van menselijk gedrag.

What could you do with data if scalability wasn't a problem? With this hands-on guide, you'll learn how Apache Cassandra handles hundreds of terabytes of data while remaining highly available across multiple data centers -- capabilities that have attracted Facebook, Twitter, and other data-intensive companies. Cassandra: The Definitive Guide provides the technical details and practical examples you need to assess this database management system and put it to work in a production environment. Author Eben Hewitt demonstrates the advantages of Cassandra's nonrelational design, and pays special attention to data modeling. If you're a developer, DBA, application architect, or manager looking to solve a database scaling issue or future-proof your application, this guide shows you how to harness Cassandra's speed and flexibility. Understand the tenets of Cassandra's column-oriented structure Learn how to write, update, and read Cassandra data Discover how to add or remove nodes from the cluster as your application requires Examine a working application that translates from a relational model to Cassandra's data model Use examples for writing clients in Java, Python, and C# Use the JMX interface to monitor a cluster's usage, memory patterns, and more Tune memory settings, data storage, and caching for better performance

This book constitutes the proceedings of the 8th International Conference on Big Data Analytics, BDA 2020, which took place during

December 15-18, 2020, in Sonapat, India. The 11 full and 3 short papers included in this volume were carefully reviewed and selected from 48 submissions; the book also contains 4 invited and 3 tutorial papers. The contributions were organized in topical sections named as follows: data science systems; data science architectures; big data analytics in healthcare; information interchange of Web data resources; and business analytics.

Business and medical professionals rely on large data sets to identify trends or other knowledge that can be gleaned from the collection of it. New technologies concentrate on data's management, but do not facilitate users' extraction of meaningful outcomes. *Pattern and Data Analysis in Healthcare Settings* investigates the approaches to shift computing from analysis on-demand to knowledge on-demand. By providing innovative tactics to apply data and pattern analysis, these practices are optimized into pragmatic sources of knowledge for healthcare professionals. This publication is an exhaustive source for policy makers, developers, business professionals, healthcare providers, and graduate students concerned with data retrieval and analysis.

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.

"Get ready to unlock the power of your data. With the fourth edition of 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."--

Apache Spark is a flexible in-memory framework that allows processing of both batch and real-time data. Its unified engine has made it quite popular for big data use cases. This book will help you to quickly get started with Apache Spark 2.0 and write efficient big data applications for a variety of use cases.

There's a lot of information about big data technologies, but splicing these technologies into an end-to-end enterprise data platform is a daunting task not widely covered. With this practical book, you'll learn how to build big data infrastructure both on-premises and in the cloud and successfully architect a modern data platform. Ideal for enterprise architects, IT managers, application architects, and data engineers, this book shows you how to overcome the many challenges that emerge during Hadoop projects. You'll explore the vast landscape of tools available in the Hadoop and big data realm in a thorough technical primer before diving into: **Infrastructure:** Look at all component layers in a modern data platform, from the server to the data center, to establish a solid foundation for data in your enterprise **Platform:** Understand aspects of deployment, operation, security, high availability, and disaster recovery, along with everything you need to know to integrate your platform with the rest of your enterprise IT **Taking Hadoop to the cloud:** Learn the important architectural aspects of running a big data platform in the cloud while maintaining enterprise security and high availability

"This book discusses the exponential growth of information size and the innovative methods for data capture, storage, sharing, and analysis for big data"--Provided by publisher.

Big data is a relative term describing a situation where the volume, velocity and variety of data exceed an organization's storage or compute capacity for accurate and timely decision making . Big data is not a single technology but a combination of old and new technologies that helps companies gain actionable insight. Therefore, big data is the capability to manage a huge volume of disparate data, at the right speed, and within the right time frame to allow real-time analysis and reaction. As we note earlier in this chapter, big data is typically broken down by three characteristics: Volume: How much data Velocity: How fast that data is processed Variety: The various types of data Although it's convenient to simplify big data into the three Vs, it can be misleading and overly simplistic. For example, you may be managing a relatively small amount of very disparate, complex data or you may be processing a huge volume of very simple data. That simple data may be all structured or all unstructured. Even more important is the fourth V: veracity. How accurate is that data in predicting business value? Do the results of a big data analysis actually make sense? Determining relevant data is key to delivering value from massive amounts of data. However, big data is defined less by volume - which is a constantly moving target - than by its ever-increasing variety, velocity, variability and complexity

Every enterprise application creates data, whether it's log messages, metrics, user activity, outgoing messages, or something else. And how to move all of this data becomes nearly as important as the data itself. If you're an application architect, developer, or production engineer new to Apache Kafka, this practical guide shows you how to use this open source streaming platform to handle real-time data feeds. Engineers from Confluent and LinkedIn who are responsible for developing Kafka explain how to deploy production Kafka clusters, write reliable event-driven microservices, and build scalable stream-processing applications with this platform. Through detailed examples, you'll learn Kafka's design principles, reliability guarantees, key APIs, and architecture details, including the replication protocol, the controller, and the storage layer. Understand publish-subscribe messaging and how it fits in the big data ecosystem. Explore Kafka producers and consumers for writing and reading messages Understand Kafka patterns and use-case requirements to ensure reliable data delivery Get best practices for building data pipelines and applications with Kafka Manage Kafka in production, and learn to perform monitoring, tuning, and maintenance tasks Learn the most critical metrics among Kafka's operational measurements Explore how Kafka's stream delivery capabilities make it a perfect source for stream processing systems

Until now, design patterns for the MapReduce framework have been scattered among various research papers, blogs, and books. This handy guide brings together a unique collection of valuable MapReduce patterns that will save you time and effort regardless of the domain, language, or development framework you're using. Each pattern is explained in context, with pitfalls and caveats clearly identified to help you avoid common design mistakes when modeling your big data architecture. This book also provides a complete overview of MapReduce that explains its origins and

implementations, and why design patterns are so important. All code examples are written for Hadoop. Summarization patterns: get a top-level view by summarizing and grouping data Filtering patterns: view data subsets such as records generated from one user Data organization patterns: reorganize data to work with other systems, or to make MapReduce analysis easier Join patterns: analyze different datasets together to discover interesting relationships Metapatterns: piece together several patterns to solve multi-stage problems, or to perform several analytics in the same job Input and output patterns: customize the way you use Hadoop to load or store data "A clear exposition of MapReduce programs for common data processing patterns--this book is indispensable for anyone using Hadoop."--Tom White, author of Hadoop: The Definitive Guide.

Get ready to unlock the power of your data. With the fourth edition of 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. Using Hadoop 2 exclusively, author Tom White presents new chapters on YARN and several Hadoop-related projects such as Parquet, Flume, Crunch, and Spark. You'll learn about recent changes to Hadoop, and explore new case studies on Hadoop's role in healthcare systems and genomics data processing. Learn fundamental components such as MapReduce, HDFS, and YARN Explore MapReduce in depth, including steps for developing applications with it Set up and maintain a Hadoop cluster running HDFS and MapReduce on YARN Learn two data formats: Avro for data serialization and Parquet for nested data Use data ingestion tools such as Flume (for streaming data) and Sqoop (for bulk data transfer) Understand how high-level data processing tools like Pig, Hive, Crunch, and Spark work with Hadoop Learn the HBase distributed database and the ZooKeeper distributed configuration service

Learn how to take full advantage of Apache Kafka, the distributed, publish-subscribe queue for handling real-time data feeds. With this comprehensive book, you'll understand how Kafka works and how it's designed. Authors Neha Narkhede, Gwen Shapira, and Todd Palino show you how to deploy production Kafka clusters; secure, tune, and monitor them; write rock-solid applications that use Kafka; and build scalable stream-processing applications. Learn how Kafka compares to other queues, and where it fits in the big data ecosystem Dive into Kafka's internal design Pick up best practices for developing applications that use Kafka Understand the best way to deploy Kafka in production monitoring, tuning, and maintenance tasks Learn how to secure a Kafka cluster Get detailed use-cases

This timely text/reference describes the development and implementation of large-scale distributed processing systems using open source tools and technologies. Comprehensive in scope, the book presents state-of-the-art material on building high performance distributed computing systems, providing practical guidance and best practices as well as

describing theoretical software frameworks. Features: describes the fundamentals of building scalable software systems for large-scale data processing in the new paradigm of high performance distributed computing; presents an overview of the Hadoop ecosystem, followed by step-by-step instruction on its installation, programming and execution; Reviews the basics of Spark, including resilient distributed datasets, and examines Hadoop streaming and working with Scalding; Provides detailed case studies on approaches to clustering, data classification and regression analysis; Explains the process of creating a working recommender system using Scalding and Spark.

Need to move a relational database application to Hadoop? This comprehensive guide introduces you to Apache Hive, Hadoop's data warehouse infrastructure. You'll quickly learn how to use Hive's SQL dialect—HiveQL—to summarize, query, and analyze large datasets stored in Hadoop's distributed filesystem. This example-driven guide shows you how to set up and configure Hive in your environment, provides a detailed overview of Hadoop and MapReduce, and demonstrates how Hive works within the Hadoop ecosystem. You'll also find real-world case studies that describe how companies have used Hive to solve unique problems involving petabytes of data. Use Hive to create, alter, and drop databases, tables, views, functions, and indexes Customize data formats and storage options, from files to external databases Load and extract data from tables—and use queries, grouping, filtering, joining, and other conventional query methods Gain best practices for creating user defined functions (UDFs) Learn Hive patterns you should use and anti-patterns you should avoid Integrate Hive with other data processing programs Use storage handlers for NoSQL databases and other datastores Learn the pros and cons of running Hive on Amazon's Elastic MapReduce

The digital age has presented an exponential growth in the amount of data available to individuals looking to draw conclusions based on given or collected information across industries. Challenges associated with the analysis, security, sharing, storage, and visualization of large and complex data sets continue to plague data scientists and analysts alike as traditional data processing applications struggle to adequately manage big data. The Handbook of Research on Big Data Storage and Visualization Techniques is a critical scholarly resource that explores big data analytics and technologies and their role in developing a broad understanding of issues pertaining to the use of big data in multidisciplinary fields. Featuring coverage on a broad range of topics, such as architecture patterns, programming systems, and computational energy, this publication is geared towards professionals, researchers, and students seeking current research and application topics on the subject.

If you're looking for a scalable storage solution to accommodate a virtually endless amount of data, this book shows you how Apache HBase can fulfill your needs. As the open source implementation of Google's BigTable architecture, HBase scales to billions of rows and millions of columns, while ensuring that write and read performance remain constant. Many

IT executives are asking pointed questions about HBase. This book provides meaningful answers, whether you're evaluating this non-relational database or planning to put it into practice right away. Discover how tight integration with Hadoop makes scalability with HBase easier. Distribute large datasets across an inexpensive cluster of commodity servers. Access HBase with native Java clients, or with gateway servers providing REST, Avro, or Thrift APIs. Get details on HBase's architecture, including the storage format, write-ahead log, background processes, and more. Integrate HBase with Hadoop's MapReduce framework for massively parallelized data processing jobs. Learn how to tune clusters, design schemas, copy tables, import bulk data, decommission nodes, and many other tasks.

This book includes the outcomes of the International Conference on Advanced Intelligent Systems for Sustainable Development (AI2SD-2018), held in Tangier, Morocco on July 12–14, 2018. Presenting the latest research in the field of computing sciences and information technology, it discusses new challenges and provides valuable insights into the field, the goal being to stimulate debate, and to promote closer interaction and interdisciplinary collaboration between researchers and practitioners. Though chiefly intended for researchers and practitioners in advanced information technology management and networking, the book will also be of interest to those engaged in emerging fields such as data science and analytics, big data, internet of things, smart networked systems, artificial intelligence, expert systems and cloud computing.

With Early Release ebooks, you get books in their earliest form—the author's raw and unedited content as he or she writes—so you can take advantage of these technologies long before the official release of these titles. You'll also receive updates when significant changes are made, new chapters as they're written, and the final ebook bundle. If you're looking for a scalable storage solution to accommodate a virtually endless amount of data, this updated edition shows you how Apache HBase can meet your needs. Modeled after Google's BigTable architecture, HBase scales to billions of rows and millions of columns, while ensuring that write and read performance remain constant. Fully revised for HBase 1.0, this second edition brings you up to speed on the new HBase client API, as well as security features and new case studies that demonstrate HBase use in the real world. Whether you just started to evaluate this non-relational database, or plan to put it into practice right away, this book has your back. Launch into basic, advanced, and administrative features of HBase's new client-facing API. Use new classes to integrate HBase with Hadoop's MapReduce framework. Explore HBase's architecture, including the storage format, write-ahead log, and background processes. Dive into advanced usage, such as extended client and server options. Learn cluster sizing, tuning, and monitoring best practices. Design schemas, copy tables, import bulk data, decommission nodes, and other tasks. Go deeper into HBase security, including Kerberos and encryption at rest.

Counsels programmers and administrators for big and small organizations on how to work with large-scale application datasets using Apache Hadoop, discussing its capacity for storing and processing large amounts of data while demonstrating best practices for building reliable and scalable distributed systems.

How can you get your data from frontend servers to Hadoop in near real time? With this complete reference guide, you'll learn Flume's rich set of features for collecting, aggregating, and writing large amounts of streaming data to the Hadoop Distributed File System (HDFS), Apache HBase, SolrCloud, Elastic Search, and other systems. Using Flume shows operations engineers how to configure, deploy, and monitor a Flume cluster, and teaches developers how to write Flume plugins and custom components for their specific use-cases. You'll learn about Flume's design and implementation, as well as various features that make it highly scalable, flexible, and reliable. Code examples and exercises are available on GitHub. Learn how Flume provides a steady rate of flow by acting as a buffer between data producers and consumers Dive into key Flume components, including sources that accept data and sinks that write and deliver it Write custom plugins to customize the way Flume receives, modifies, formats, and writes data Explore APIs for sending data to Flume agents from your own applications Plan and deploy Flume in a scalable and flexible way—and monitor your cluster once it's running

This multi-contributed handbook focuses on the latest workings of IoT (internet of Things) and Big Data. As the resources are limited, it's the endeavor of the authors to support and bring the information into one resource. The book is divided into 4 sections that covers IoT and technologies, the future of Big Data, algorithms, and case studies showing IoT and Big Data in various fields such as health care, manufacturing and automation. Features Focuses on the latest workings of IoT and Big Data Discusses the emerging role of technologies and the fast-growing market of Big Data Covers the movement toward automation with hardware, software, and sensors, and trying to save on energy resources Offers the latest technology on IoT Presents the future horizons on Big Data

This proceedings volume brings together peer-reviewed papers presented at the International Conference on Information Technology and Computer Application Engineering, held 10-11 December 2014, in Hong Kong, China. Specific topics under consideration include Computational Intelligence, Computer Science and its Applications, Intelligent Information Processing and Knowledge Engineering, Intelligent Networks and Instruments, Multimedia Signal Processing and Analysis, Intelligent Computer-Aided Design Systems and other related topics. This book provides readers a state-of-the-art survey of recent innovations and research worldwide in Information Technology and Computer Application Engineering, in so-doing furthering the development and growth of these research fields, strengthening international academic cooperation and communication, and promoting the fruitful exchange of research ideas. This volume will be of interest to professionals and academics alike, serving as a broad overview of the latest advances in the dynamic field of Information Technology and Computer Application Engineering.

Hadoop: The Definitive Guide"O'Reilly Media, Inc."

Until recently, Hadoop deployments existed on hardware owned and run by organizations. Now, of course, you can acquire the computing resources and network connectivity to run Hadoop clusters in the cloud. But there's a lot more to deploying Hadoop to the public cloud than simply renting machines. This hands-on guide shows developers and systems administrators familiar with Hadoop how to install, use, and manage cloud-born clusters efficiently. You'll learn how to architect clusters that work with cloud-

provider features—not just to avoid pitfalls, but also to take full advantage of these services. You'll also compare the Amazon, Google, and Microsoft clouds, and learn how to set up clusters in each of them. Learn how Hadoop clusters run in the cloud, the problems they can help you solve, and their potential drawbacks Examine the common concepts of cloud providers, including compute capabilities, networking and security, and storage Build a functional Hadoop cluster on cloud infrastructure, and learn what the major providers require Explore use cases for high availability, relational data with Hive, and complex analytics with Spark Get patterns and practices for running cloud clusters, from designing for price and security to dealing with maintenance

[Copyright: 97e640db9209d0eda89cfc1a1a98798d](#)