

# Learn Apache Tika Java Technologies

Over 100 practical recipes to make Apache Solr faster, more reliable and return better results.

The inside story of how America's enemies launched a cyber war against us-and how we've learned to fight back With each passing year, the internet-linked attacks on America's interests have grown in both frequency and severity.

Overmatched by our military, countries like North Korea, China, Iran, and Russia have found us vulnerable in cyberspace. The "Code War" is upon us. In this dramatic book, former Assistant Attorney General John P. Carlin takes readers to the front lines of a global but little-understood fight as the Justice Department and the FBI chases down hackers, online terrorist recruiters, and spies. Today, as our entire economy goes digital, from banking to manufacturing to transportation, the potential targets for our enemies multiply. This firsthand account is both a remarkable untold story and a warning of dangers yet to come.

Summary Deep Learning for Search teaches you how to improve the effectiveness of your search by implementing neural network-based techniques. By the time you're finished with the book, you'll be ready to build amazing search engines that deliver the results your users need and that get better as time goes on! Foreword by Chris Mattmann.

Purchase of the print book includes a free eBook in PDF, Kindle, and ePub formats from Manning Publications. About the Technology Deep learning handles the toughest search challenges, including imprecise search terms, badly indexed data, and retrieving images with minimal metadata. And with modern tools like DL4J and TensorFlow, you can apply powerful DL techniques without a deep background in data

science or natural language processing (NLP). This book will show you how. About the Book Deep Learning for Search teaches you to improve your search results with neural networks. You'll review how DL relates to search basics like indexing and ranking. Then, you'll walk through in-depth examples to upgrade your search with DL techniques using Apache Lucene and Deeplearning4j. As the book progresses, you'll explore advanced topics like searching through images, translating user queries, and designing search engines that improve as they learn! What's inside Accurate and relevant rankings Searching across languages Content-based image search Search with recommendations About the Reader For developers comfortable with Java or a similar language and search basics. No experience with deep learning or NLP needed. About the Author Tommaso Teofili is a software engineer with a passion for open source and machine learning. As a member of the Apache Software Foundation, he contributes to a number of open source projects, ranging from topics like information retrieval (such as Lucene and Solr) to natural language processing and machine translation (including OpenNLP, Joshua, and UIMA). He currently works at Adobe, developing search and indexing infrastructure components, and researching the areas of natural language processing, information retrieval, and deep learning. He has presented search and machine learning talks at conferences including BerlinBuzzwords, International Conference on Computational Science, ApacheCon, EclipseCon, and others. You can find him on Twitter at @tteofili. Table of Contents PART 1 - SEARCH MEETS DEEP LEARNING Neural search Generating synonyms PART 2 - THROWING NEURAL NETS AT A SEARCH ENGINE From plain retrieval to text generation More-sensitive query suggestions Ranking search results with word embeddings Document embeddings for rankings and recommendations PART 3 - ONE STEP

BEYOND Searching across languages Content-based image search A peek at performance

Hadoop in Action teaches readers how to use Hadoop and write MapReduce programs. The intended readers are programmers, architects, and project managers who have to process large amounts of data offline. Hadoop in Action will lead the reader from obtaining a copy of Hadoop to setting it up in a cluster and writing data analytic programs. The book begins by making the basic idea of Hadoop and MapReduce easier to grasp by applying the default Hadoop installation to a few easy-to-follow tasks, such as analyzing changes in word frequency across a body of documents. The book continues through the basic concepts of MapReduce applications developed using Hadoop, including a close look at framework components, use of Hadoop for a variety of data analysis tasks, and numerous examples of Hadoop in action. Hadoop in Action will explain how to use Hadoop and present design patterns and practices of programming MapReduce. MapReduce is a complex idea both conceptually and in its implementation, and Hadoop users are challenged to learn all the knobs and levers for running Hadoop. This book takes you beyond the mechanics of running Hadoop, teaching you to write meaningful programs in a MapReduce framework. This book assumes the reader will have a basic familiarity with Java, as most code examples will be written in Java. Familiarity with basic statistical concepts (e.g. histogram, correlation) will help the reader appreciate the more advanced data processing examples. Purchase of the print book comes with an offer of a free PDF, ePub, and Kindle eBook from Manning. Also available is all code from the book.

This two-volume book contains research work presented at the First International Conference on Data Engineering and Communication Technology (ICDECT) held during March

10–11, 2016 at Lavasa, Pune, Maharashtra, India. The book discusses recent research technologies and applications in the field of Computer Science, Electrical and Electronics Engineering. The aim of the Proceedings is to provide cutting-edge developments taking place in the field data engineering and communication technologies which will assist the researchers and practitioners from both academia as well as industry to advance their field of study.

Ongoing advancements in modern technology have led to significant developments in artificial intelligence. With the numerous applications available, it becomes imperative to conduct research and make further progress in this field. *Artificial Intelligence: Concepts, Methodologies, Tools, and Applications* provides a comprehensive overview of the latest breakthroughs and recent progress in artificial intelligence. Highlighting relevant technologies, uses, and techniques across various industries and settings, this publication is a pivotal reference source for researchers, professionals, academics, upper-level students, and practitioners interested in emerging perspectives in the field of artificial intelligence. *Dig deep into the data with a hands-on guide to machine learning with updated examples and more!* *Machine Learning: Hands-On for Developers and Technical Professionals* provides hands-on instruction and fully-coded working examples for the most common machine learning techniques used by developers and technical professionals. The book contains a breakdown of each ML variant, explaining how it works and how it is used within certain industries, allowing readers to incorporate the presented techniques into their own work as they follow along. A core tenant of machine learning is a strong focus on data preparation, and a full exploration of the various types of learning algorithms illustrates how the proper tools can help any developer extract information and insights from existing data. The book

includes a full complement of Instructor's Materials to facilitate use in the classroom, making this resource useful for students and as a professional reference. At its core, machine learning is a mathematical, algorithm-based technology that forms the basis of historical data mining and modern big data science. Scientific analysis of big data requires a working knowledge of machine learning, which forms predictions based on known properties learned from training data. Machine Learning is an accessible, comprehensive guide for the non-mathematician, providing clear guidance that allows readers to: Learn the languages of machine learning including Hadoop, Mahout, and Weka Understand decision trees, Bayesian networks, and artificial neural networks Implement Association Rule, Real Time, and Batch learning Develop a strategic plan for safe, effective, and efficient machine learning By learning to construct a system that can learn from data, readers can increase their utility across industries. Machine learning sits at the core of deep dive data analysis and visualization, which is increasingly in demand as companies discover the goldmine hiding in their existing data. For the tech professional involved in data science, Machine Learning: Hands-On for Developers and Technical Professionals provides the skills and techniques required to dig deeper.

Summary Taming Text, winner of the 2013 Jolt Awards for Productivity, is a hands-on, example-driven guide to working with unstructured text in the context of real-world applications. This book explores how to automatically organize text using approaches such as full-text search, proper name recognition, clustering, tagging, information extraction, and summarization. The book guides you through examples illustrating each of these topics, as well as the foundations upon which they are built. About this Book There is so much text in our lives, we are practically drowning in it. Fortunately,

there are innovative tools and techniques for managing unstructured information that can throw the smart developer a much-needed lifeline. You'll find them in this book. *Taming Text* is a practical, example-driven guide to working with text in real applications. This book introduces you to useful techniques like full-text search, proper name recognition, clustering, tagging, information extraction, and summarization. You'll explore real use cases as you systematically absorb the foundations upon which they are built. Written in a clear and concise style, this book avoids jargon, explaining the subject in terms you can understand without a background in statistics or natural language processing. Examples are in Java, but the concepts can be applied in any language. Written for Java developers, the book requires no prior knowledge of GWT. Purchase of the print book comes with an offer of a free PDF, ePub, and Kindle eBook from Manning. Also available is all code from the book. Winner of 2013 Jolt Awards: The Best Books—one of five notable books every serious programmer should read.

**What's Inside**

- When to use text-taming techniques
- Important open-source libraries like Solr and Mahout
- How to build text-processing applications
- About the Authors

**Grant Ingersoll** is an engineer, speaker, and trainer, a Lucene committer, and a cofounder of the Mahout machine-learning project. **Thomas Morton** is the primary developer of OpenNLP and Maximum Entropy. **Drew Farris** is a technology consultant, software developer, and contributor to Mahout, Lucene, and Solr.

"Takes the mystery out of very complex processes."—From the Foreword by Liz Liddy, Dean, iSchool, Syracuse University

**Table of Contents**

- Getting started taming text
- Foundations of taming text
- Searching
- Fuzzy string matching
- Identifying people, places, and things
- Clustering text
- Classification, categorization, and tagging
- Building an example question answering system
- Untamed text: exploring the next frontier

Tika in Action Simon and Schuster

"XQuery Kick Start" delivers a concise introduction to the XQuery standard, and useful implementation advice for developers needing to put it into practice. The book starts by explaining the role of XQuery in the XML family of specifications, and its relationship with XPath. The authors then explain the specification in detail, describing the semantics and data model, before moving to examples using XQuery to manipulate XML databases and document storage systems. Later chapters discuss Java implementations of XQuery and development tools that facilitate the development of Web sites with XQuery. This book is up to date with the latest XQuery specifications, and includes coverage of new features for extending the XQuery language.

This book is for developers who want to learn how to get the most out of Solr in their applications, whether you are new to the field, have used Solr but don't know everything, or simply want a good reference. It would be helpful to have some familiarity with basic programming concepts, but no prior experience is required.

A practical guide to implementing your enterprise data lake using Lambda Architecture as the base About This Book Build a full-fledged data lake for your organization with popular big data technologies using the Lambda architecture as the base Delve into the big data technologies required to meet modern day business strategies A highly practical guide to implementing enterprise data lakes with lots of examples and real-world use-cases Who This Book Is For Java developers and architects who would like to implement a data lake

for their enterprise will find this book useful. If you want to get hands-on experience with the Lambda Architecture and big data technologies by implementing a practical solution using these technologies, this book will also help you. What You Will Learn Build an enterprise-level data lake using the relevant big data technologies Understand the core of the Lambda architecture and how to apply it in an enterprise Learn the technical details around Sqoop and its functionalities Integrate Kafka with Hadoop components to acquire enterprise data Use flume with streaming technologies for stream-based processing Understand stream-based processing with reference to Apache Spark Streaming Incorporate Hadoop components and know the advantages they provide for enterprise data lakes Build fast, streaming, and high-performance applications using Elasticsearch Make your data ingestion process consistent across various data formats with configurability Process your data to derive intelligence using machine learning algorithms In Detail The term "Data Lake" has recently emerged as a prominent term in the big data industry. Data scientists can make use of it in deriving meaningful insights that can be used by businesses to redefine or transform the way they operate. Lambda architecture is also emerging as one of the very eminent patterns in the big data landscape, as it not only helps to derive useful information from historical data but also correlates real-time data to enable business to take critical decisions. This book tries to bring these two important aspects — data lake and lambda architecture—together. This book is divided into

three main sections. The first introduces you to the concept of data lakes, the importance of data lakes in enterprises, and getting you up-to-speed with the Lambda architecture. The second section delves into the principal components of building a data lake using the Lambda architecture. It introduces you to popular big data technologies such as Apache Hadoop, Spark, Sqoop, Flume, and ElasticSearch. The third section is a highly practical demonstration of putting it all together, and shows you how an enterprise data lake can be implemented, along with several real-world use-cases. It also shows you how other peripheral components can be added to the lake to make it more efficient. By the end of this book, you will be able to choose the right big data technologies using the lambda architectural patterns to build your enterprise data lake. Style and approach The book takes a pragmatic approach, showing ways to leverage big data technologies and lambda architecture to build an enterprise-level data lake.

Explore various approaches to organize and extract useful text from unstructured data using Java Key Features Use deep learning and NLP techniques in Java to discover hidden insights in text Work with popular Java libraries such as CoreNLP, OpenNLP, and Mallet Explore machine translation, identifying parts of speech, and topic modeling Book Description Natural Language Processing (NLP) allows you to take any sentence and identify patterns, special names, company names, and more. The second edition of Natural Language Processing with Java teaches you how to perform language analysis with the help of Java libraries, while

constantly gaining insights from the outcomes. You'll start by understanding how NLP and its various concepts work. Having got to grips with the basics, you'll explore important tools and libraries in Java for NLP, such as CoreNLP, OpenNLP, Neuroph, and Mallet. You'll then start performing NLP on different inputs and tasks, such as tokenization, model training, parts-of-speech and parsing trees. You'll learn about statistical machine translation, summarization, dialog systems, complex searches, supervised and unsupervised NLP, and more. By the end of this book, you'll have learned more about NLP, neural networks, and various other trained models in Java for enhancing the performance of NLP applications. What you will learn

- Understand basic NLP tasks and how they relate to one another
- Discover and use the available tokenization engines
- Apply search techniques to find people, as well as things, within a document
- Construct solutions to identify parts of speech within sentences
- Use parsers to extract relationships between elements of a document
- Identify topics in a set of documents
- Explore topic modeling from a document

Who this book is for

Natural Language Processing with Java is for you if you are a data analyst, data scientist, or machine learning engineer who wants to extract information from a language using Java. Knowledge of Java programming is needed, while a basic understanding of statistics will be useful but not mandatory.

Summary

Introducing Data Science teaches you how to accomplish the fundamental tasks that occupy data scientists. Using the Python language and common

Python libraries, you'll experience firsthand the challenges of dealing with data at scale and gain a solid foundation in data science. Purchase of the print book includes a free eBook in PDF, Kindle, and ePub formats from Manning Publications. About the Technology Many companies need developers with data science skills to work on projects ranging from social media marketing to machine learning. Discovering what you need to learn to begin a career as a data scientist can seem bewildering. This book is designed to help you get started. About the Book *Introducing Data Science* *Introducing Data Science* explains vital data science concepts and teaches you how to accomplish the fundamental tasks that occupy data scientists. You'll explore data visualization, graph databases, the use of NoSQL, and the data science process. You'll use the Python language and common Python libraries as you experience firsthand the challenges of dealing with data at scale. Discover how Python allows you to gain insights from data sets so big that they need to be stored on multiple machines, or from data moving so quickly that no single machine can handle it. This book gives you hands-on experience with the most popular Python data science libraries, Scikit-learn and StatsModels. After reading this book, you'll have the solid foundation you need to start a career in data science. What's Inside Handling large data Introduction to machine learning Using Python to work with data Writing data science algorithms About the Reader This book assumes you're comfortable reading code in Python or a similar language, such as C, Ruby, or JavaScript. No prior experience with data science is

required. About the Authors Davy Cielen, Arno D. B. Meysman, and Mohamed Ali are the founders and managing partners of Optimately and Maiton, where they focus on developing data science projects and solutions in various sectors. Table of Contents Data science in a big data world The data science process Machine learning Handling large data on a single computer First steps in big data Join the NoSQL movement The rise of graph databases Text mining and text analytics Data visualization to the end user

This book is a step-by-step guide for readers who would like to learn how to build complete enterprise search solutions, with ample real-world examples and case studies. If you are a developer, designer, or architect who would like to build enterprise search solutions for your customers or organization, but have no prior knowledge of Apache Solr/Lucene technologies, this is the book for you.

Summary Solr in Action is a comprehensive guide to implementing scalable search using Apache Solr. This clearly written book walks you through well-documented examples ranging from basic keyword searching to scaling a system for billions of documents and queries. It will give you a deep understanding of how to implement core Solr capabilities. About the Book Whether you're handling big (or small) data, managing documents, or building a website, it is important to be able to quickly search through your content and discover meaning in it. Apache Solr is your tool: a ready-to-deploy, Lucene-based, open source, full-text search engine. Solr can scale across many servers to enable real-time queries

and data analytics across billions of documents. Solr in Action teaches you to implement scalable search using Apache Solr. This easy-to-read guide balances conceptual discussions with practical examples to show you how to implement all of Solr's core capabilities. You'll master topics like text analysis, faceted search, hit highlighting, result grouping, query suggestions, multilingual search, advanced geospatial and data operations, and relevancy tuning. This book assumes basic knowledge of Java and standard database technology. No prior knowledge of Solr or Lucene is required. Purchase of the print book includes a free eBook in PDF, Kindle, and ePub formats from Manning Publications.

What's Inside

- How to scale Solr for big data
- Rich real-world examples
- Solr as a NoSQL data store
- Advanced multilingual, data, and relevancy tricks
- Coverage of versions through Solr 4.7

About the Authors

Trey Grainger is a director of engineering at CareerBuilder. Timothy Potter is a senior member of the engineering team at LucidWorks. The authors work on the scalability and reliability of Solr, as well as on recommendation engine and big data analytics technologies.

Table of Contents

PART 1 MEET SOLR

- Introduction to Solr
- Getting to know Solr
- Key Solr concepts
- Configuring Solr
- Indexing
- Text analysis

PART 2 CORE SOLR CAPABILITIES

- Performing queries and handling results
- Faceted search
- Hit highlighting
- Query suggestions
- Result grouping/field collapsing
- Taking Solr to production

PART 3 TAKING SOLR TO THE NEXT LEVEL

- SolrCloud
- Multilingual search
- Complex query operations
- Mastering relevancy

Accelerate your enterprise search engine and bring relevancy in your search analytics Key Features A practical guide in building expertise with Indexing, Faceting, Clustering and Pagination Master the management and administration of Enterprise Search Applications and services seamlessly Handle multiple data inputs such as JSON, xml, pdf, doc, xls,ppt, csv and much more. Book Description Apache Solr is the only standalone enterprise search server with a REST-like application interface. providing highly scalable, distributed search and index replication for many of the world's largest internet sites. To begin with, you would be introduced to how you perform full text search, multiple filter search, perform dynamic clustering and so on helping you to brush up the basics of Apache Solr. You will also explore the new features and advanced options released in Apache Solr 7.x which will get you numerous performance aspects and making data investigation simpler, easier and powerful. You will learn to build complex queries, extensive filters and how are they compiled in your system to bring relevance in your search tools. You will learn to carry out Solr scoring, elements affecting the document score and how you can optimize or tune the score for the application at hand. You will learn to extract features of documents, writing complex queries in re-ranking the documents. You will also learn advanced options helping you to know what content is indexed and how the extracted content is indexed. Throughout the book, you would go through complex problems with solutions along with varied approaches to tackle your business needs. By the end of

this book, you will gain advanced proficiency to build out-of-box smart search solutions for your enterprise demands. What you will learn Design schema using schema API to access data in the database Advance querying and fine-tuning techniques for better performance Get to grips with indexing using Client API Set up a fault tolerant and highly available server with newer distributed capabilities, SolrCloud Explore Apache Tika to upload data with Solr Cell Understand different data operations that can be done while indexing Master advanced querying through Velocity Search UI, faceting and Query Re-ranking, pagination and spatial search Learn to use JavaScript, Python, SolrJ and Ruby for interacting with Solr Who this book is for The book would rightly appeal to developers, software engineers, data engineers and database architects who are building or seeking to build enterprise-wide effective search engines for business intelligence. Prior experience of Apache Solr or Java programming is must to take the best of this book.

Search is everywhere, yet it is one of the most misunderstood functionalities of the IT industry. In Apache Solr Succinctly, author Xavier Morera guides you through the basics of this highly popular enterprise search tool. You'll learn how to set up an index and how to make it searchable, then query it with a simple enterprise search. Explanations for precision and recall are also included to help you ensure that relevant, accurate results have been returned. Custom UIs using SolrItas and SolrNet are also covered.

This book constitutes the thoroughly refereed post-

workshop proceedings of the 4th International Symposium, SETE 2019, held in conjunction with ICWL 2019, in Magdeburg, Germany, in September 2019. The 10 full and 6 short papers presented together with 24 papers from 5 workshops were carefully reviewed and selected from 34 submissions. The papers cover the latest findings in various areas, such as: virtual reality and game-based learning; learning analytics; K-12 education; language learning; design, model and implementation of e-learning platforms and tools; digitalization and industry 4.0; pedagogical issues, practice and experience sharing.

This report improves the evidence base on the role of Data Driven Innovation for promoting growth and well-being, and provide policy guidance on how to maximise the benefits of DDI and mitigate the associated economic and societal risks.

Get up to speed on the nuances of NoSQL databases and what they mean for your organization This easy to read guide to NoSQL databases provides the type of no-nonsense overview and analysis that you need to learn, including what NoSQL is and which database is right for you. Featuring specific evaluation criteria for NoSQL databases, along with a look into the pros and cons of the most popular options, NoSQL For Dummies provides the fastest and easiest way to dive into the details of this incredible technology. You'll gain an understanding of how to use NoSQL databases for mission-critical enterprise architectures and projects, and real-world examples reinforce the primary points to create an action-oriented resource for IT pros. If you're planning a big

data project or platform, you probably already know you need to select a NoSQL database to complete your architecture. But with options flooding the market and updates and add-ons coming at a rapid pace, determining what you require now, and in the future, can be a tall task. This is where NoSQL For Dummies comes in! Learn the basic tenets of NoSQL databases and why they have come to the forefront as data has outpaced the capabilities of relational databases Discover major players among NoSQL databases, including Cassandra, MongoDB, MarkLogic, Neo4J, and others Get an in-depth look at the benefits and disadvantages of the wide variety of NoSQL database options Explore the needs of your organization as they relate to the capabilities of specific NoSQL databases Big data and Hadoop get all the attention, but when it comes down to it, NoSQL databases are the engines that power many big data analytics initiatives. With NoSQL For Dummies, you'll go beyond relational databases to ramp up your enterprise's data architecture in no time.

Summary Tika in Action is a hands-on guide to content mining with Apache Tika. The book's many examples and case studies offer real-world experience from domains ranging from search engines to digital asset management and scientific data processing. About the Technology Tika is an Apache toolkit that has built into it everything you and your app need to know about file formats. Using Tika, your applications can discover and extract content from digital documents in almost any format, including exotic ones. About this Book Tika in Action is the ultimate guide to content mining using

Apache Tika. You'll learn how to pull usable information from otherwise inaccessible sources, including internet media and file archives. This example-rich book teaches you to build and extend applications based on real-world experience with search engines, digital asset management, and scientific data processing. In addition to architectural overviews, you'll find detailed chapters on features like metadata extraction, automatic language detection, and custom parser development. This book is written for developers who are new to both Scala and Lift and covers just enough Scala to get you started.

Purchase of the print book comes with an offer of a free PDF, ePub, and Kindle eBook from Manning. Also available is all code from the book. What's Inside Crack MS Word, PDF, HTML, and ZIP Integrate with search engines, CMS, and other data sources Learn through experimentation Many examples This book requires no previous knowledge of Tika or text mining techniques. It assumes a working knowledge of Java.

=====?

Table of Contents PART 1 GETTING STARTED The case for the digital Babel fish Getting started with Tika The information landscape PART 2 TIKA IN DETAIL Document type detection Content extraction Understanding metadata Language detection What's in a file? PART 3 INTEGRATION AND ADVANCED USE The big picture Tika and the Lucene search stack Extending Tika PART 4 CASE STUDIES Powering NASA science data systems Content management with Apache Jackrabbit Curating cancer research data with Tika The classic search engine example

This book constitutes the refereed proceedings of the 4th IFIP WG 5.5/SOCOLNET Doctoral Conference on Computing, Electrical and Industrial Systems, DoCEIS 2013, held in Costa de Caparica, Portugal, in April 2013. The 69 revised full papers were carefully reviewed and selected from numerous submissions. They cover a wide spectrum of topics ranging from collaborative enterprise networks to microelectronics. The papers are organized in the following topical sections: collaborative enterprise networks; service orientation; intelligent computational systems; computational systems; computational systems applications; perceptual systems; robotics and manufacturing; embedded systems and Petri nets; control and decision; integration of power electronics systems with ICT; energy generation; energy distribution; energy transformation; optimization techniques in energy; telecommunications; electronics: devices design; electronics: amplifiers; electronics: RF applications; and electronics: applications.

Designing and writing a real-time streaming publication with Apache Apex About This Book Get a clear, practical approach to real-time data processing Program Apache Apex streaming applications This book shows you Apex integration with the open source Big Data ecosystem Who This Book Is For This book assumes knowledge of application development with Java and familiarity with distributed systems. Familiarity with other real-time streaming frameworks is not required, but some practical experience with other big data processing utilities might be helpful. What You Will Learn Put together a functioning Apex application from scratch Scale an Apex

application and configure it for optimal performance  
Understand how to deal with failures via the fault tolerance features of the platform Use Apex via other frameworks such as Beam Understand the DevOps implications of deploying Apex In Detail Apache Apex is a next-generation stream processing framework designed to operate on data at large scale, with minimum latency, maximum reliability, and strict correctness guarantees. Half of the book consists of Apex applications, showing you key aspects of data processing pipelines such as connectors for sources and sinks, and common data transformations. The other half of the book is evenly split into explaining the Apex framework, and tuning, testing, and scaling Apex applications. Much of our economic world depends on growing streams of data, such as social media feeds, financial records, data from mobile devices, sensors and machines (the Internet of Things - IoT). The projects in the book show how to process such streams to gain valuable, timely, and actionable insights. Traditional use cases, such as ETL, that currently consume a significant chunk of data engineering resources are also covered. The final chapter shows you future possibilities emerging in the streaming space, and how Apache Apex can contribute to it. Style and approach This book is divided into two major parts: first it explains what Apex is, what its relevant parts are, and how to write well-built Apex applications. The second part is entirely application-driven, walking you through Apex applications of increasing complexity.

Build an enterprise search engine using Apache Solr: index

and search documents; ingest data from varied sources; apply various text processing techniques; utilize different search capabilities; and customize Solr to retrieve the desired results. Apache Solr: A Practical Approach to Enterprise Search explains each essential concept--backed by practical and industry examples--to help you attain expert-level knowledge. The book, which assumes a basic knowledge of Java, starts with an introduction to Solr, followed by steps to setting it up, indexing your first set of documents, and searching them. It then introduces you to information retrieval and its implementation in Apache Solr; this will help you understand your search problem, decide the approach to build an effective solution, and use various metrics to evaluate the results. The book next covers the schema design and techniques to build a text analysis chain for cleansing, normalizing and enriching your documents and addressing different types of search queries. It describes various popular matching techniques which are generally applied to improve the precision and recall of searches. You will learn the end-to-end process of data ingestion from varied sources, metadata extraction, pre-processing and transformation of content, various search components, query parsers and other advanced search capabilities. After covering out-of-the-box features, Solr expert Dikshant Shahi dives into ways you can customize Solr for your business and its specific requirements, along with ways to plug in your own components. Most important, you will learn about implementations for Solr scoring, factors affecting the document score, and tuning the score for the application at hand. The book explains why textual scoring is not sufficient for practical ranking of documents and ways to integrate real-world factors for contributing to the document ranking. You'll see how to influence user experience by providing suggestions and recommendations. You'll also see

integration of Solr with important related technologies such as OpenNLP and Tika. Additionally, you will learn about scaling Solr using SolrCloud. This book concludes with coverage of semantic search capabilities, which is crucial for taking the search experience to the next level. By the end of Apache Solr, you will be proficient in designing and developing your search engine.

Enhance your Solr indexing experience with advanced techniques and the built-in functionalities available in Apache Solr About This Book Learn about distributed indexing and real-time optimization to change index data on fly Index data from various sources and web crawlers using built-in analyzers and tokenizers This step-by-step guide is packed with real-life examples on indexing data Who This Book Is For This book is for developers who want to increase their experience of indexing in Solr by learning about the various index handlers, analyzers, and methods available in Solr. Beginner level Solr development skills are expected. What You Will Learn Get to know the basic features of Solr indexing and the analyzers/tokenizers available Index XML/JSON data in Solr using the HTTP Post tool and CURL command Work with Data Import Handler to index data from a database Use Apache Tika with Solr to index word documents, PDFs, and much more Utilize Apache Nutch and Solr integration to index crawled data from web pages Update indexes in real-time data feeds Discover techniques to index multi-language and distributed data in Solr Combine the various indexing techniques into a real-life working example of an online shopping web application In Detail Apache Solr is a widely used, open source enterprise search server that delivers powerful indexing and searching features. These features help fetch relevant information from various sources and documentation. Solr also combines with other open source tools such as Apache Tika and Apache Nutch to

provide more powerful features. This fast-paced guide starts by helping you set up Solr and get acquainted with its basic building blocks, to give you a better understanding of Solr indexing. You'll quickly move on to indexing text and boosting the indexing time. Next, you'll focus on basic indexing techniques, various index handlers designed to modify documents, and indexing a structured data source through Data Import Handler. Moving on, you will learn techniques to perform real-time indexing and atomic updates, as well as more advanced indexing techniques such as de-duplication. Later on, we'll help you set up a cluster of Solr servers that combine fault tolerance and high availability. You will also gain insights into working scenarios of different aspects of Solr and how to use Solr with e-commerce data. By the end of the book, you will be competent and confident working with indexing and will have a good knowledge base to efficiently program elements. Style and approach This fast-paced guide is packed with examples that are written in an easy-to-follow style, and are accompanied by detailed explanation. Working examples are included to help you get better results for your applications.

This book includes selected papers from the 2nd International Conference on Machine Learning and Information Processing (ICMLIP 2020), held at Vardhaman College of Engineering, Jawaharlal Nehru Technological University (JNTU), Hyderabad, India, from November 28 to 29, 2020. It presents the latest developments and technical solutions in the areas of advanced computing and data sciences, covering machine learning, artificial intelligence, human-computer interaction, IoT, deep learning, image processing and pattern recognition, and signal and speech processing.

Use digital experience platforms (DXP) to improve your development productivity and release timelines. Leverage the pre-integrated feature sets of DXPs in your organization's

digital transformation journey to quickly develop a personalized, secure, and robust enterprise platform. In this book the authors examine various features of DXPs and provide rich insights into building each layer in a digital platform. Proven best practices are presented with examples for designing and building layers. A special focus is provided on security and quality attributes needed for business-critical enterprise applications. The authors cover modern and emerging digital trends such as Blockchain, IoT, containers, chatbots, artificial intelligence, and more. The book is divided into five parts related to requirements/design, development, security, infrastructure, and case study. The authors employ proven real-world methods, best practices, and security and integration techniques derived from their rich experience. An elaborate digital transformation case study for a banking application is included. What You'll Learn Develop a digital experience platform from end to end Understand best practices and proven methods for designing overall architecture, user interface and integration components, security, and infrastructure Study real-world cases, including an elaborate digital transformation building an enterprise platform for a banking application Know the open source tools and technology frameworks that can be used to build DXPs Who This Book Is For Web developers, full stack developers, digital enthusiasts, digital project managers, and architects This volume includes 73 papers presented at ICTIS 2017: Second International Conference on Information and Communication Technology for Intelligent Systems. The conference was held on 25th and 26th March 2017, in Ahmedabad, India and organized jointly by the Associated Chambers of Commerce and Industry of India (ASSOCHAM) Gujarat Chapter, the G R Foundation, the Association of Computer Machinery, Ahmedabad Chapter and supported by the Computer Society of India Division IV – Communication

and Division V – Education and Research. The papers featured mainly focus on information and communications technology (ICT) and its applications in intelligent computing, cloud storage, data mining and software analysis. The fundamentals of various data analytics and algorithms discussed are useful to researchers in the field.

If you are a system or application developer interested in learning how to solve practical problems using the Hadoop framework, then this book is ideal for you. You are expected to be familiar with the Unix/Linux command-line interface and have some experience with the Java programming language. Familiarity with Hadoop would be a plus.

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.

Industrial revolutions have impacted both, manufacturing and service. From the steam engine to digital automated production, the industrial revolutions have conducted significant changes in operations and supply chain management (SCM) processes. Swift changes in manufacturing and service systems have led to phenomenal improvements in productivity. The fast-paced environment brings new challenges and opportunities for the companies that are associated with the adaptation to the new concepts such as Internet of Things (IoT) and Cyber Physical Systems, artificial intelligence (AI), robotics, cyber security, data analytics, block chain and cloud technology. These emerging technologies facilitated and expedited the birth of Logistics 4.0. Industrial Revolution 4.0 initiatives in SCM has attracted stakeholders' attentions due to its ability to empower using a set of technologies together that helps to execute more efficient production and distribution systems. This initiative has been called Logistics 4.0 of the fourth Industrial Revolution in SCM due to its high potential. Connecting entities, machines, physical items and enterprise resources to each other by using sensors, devices and the internet along the supply chains are the main attributes of Logistics 4.0. IoT enables customers to make more suitable and valuable decisions due to the data-driven structure of the Industry 4.0 paradigm. Besides that, the system's ability of gathering and analyzing information about the environment at any given time and adapting itself to the rapid changes add significant

value to the SCM processes. In this peer-reviewed book, experts from all over the world, in the field present a conceptual framework for Logistics 4.0 and provide examples for usage of Industry 4.0 tools in SCM. This book is a work that will be beneficial for both practitioners and students and academicians, as it covers the theoretical framework, on the one hand, and includes examples of practice and real world. Enterprise Integration Patterns provides an invaluable catalog of sixty-five patterns, with real-world solutions that demonstrate the formidable of messaging and help you to design effective messaging solutions for your enterprise. The authors also include examples covering a variety of different integration technologies, such as JMS, MSMQ, TIBCO ActiveEnterprise, Microsoft BizTalk, SOAP, and XSL. A case study describing a bond trading system illustrates the patterns in practice, and the book offers a look at emerging standards, as well as insights into what the future of enterprise integration might hold. This book provides a consistent vocabulary and visual notation framework to describe large-scale integration solutions across many technologies. It also explores in detail the advantages and limitations of asynchronous messaging architectures. The authors present practical advice on designing code that connects an application to a messaging system, and provide extensive information to help you determine when to send a message, how to route it to the proper destination, and how to monitor the health of a messaging system. If you want to know how to manage, monitor, and maintain a messaging system once it is in use, get this book.

When Lucene first hit the scene five years ago, it was nothing short of amazing. By using this open-source, highly scalable, super-fast search engine, developers could integrate search into applications quickly and efficiently. A lot has changed since then—search has grown from a "nice-to-have" feature into

an indispensable part of most enterprise applications. Lucene now powers search in diverse companies including Akamai, Netflix, LinkedIn, Technorati, HotJobs, Epiphany, FedEx, Mayo Clinic, MIT, New Scientist Magazine, and many others. Some things remain the same, though. Lucene still delivers high-performance search features in a disarmingly easy-to-use API. Due to its vibrant and diverse open-source community of developers and users, Lucene is relentlessly improving, with evolutions to APIs, significant new features such as payloads, and a huge increase (as much as 8x) in indexing speed with Lucene 2.3. And with clear writing, reusable examples, and unmatched advice on best practices, Lucene in Action, Second Edition is still the definitive guide to developing with Lucene. Purchase of the print book comes with an offer of a free PDF, ePub, and Kindle eBook from Manning. Also available is all code from the book.

Summary CMIS and Apache Chemistry in Action is a comprehensive guide to the CMIS standard and related ECM concepts, written by the authors of the standard. In it, you'll tackle hands-on examples for building applications on CMIS repositories from both the client and the server sides. You'll learn how to create new content-centric applications that install and run in any CMIS-compliant repository. About The Technology Content Management Interoperability Services (CMIS) is an OASIS standard for accessing content management systems. It specifies a vendor- and language-neutral way to interact with any compliant content repository. Apache Chemistry provides complete reference implementations of the CMIS standard with robust APIs for developers writing tools, applications, and servers. About This Book CMIS and Apache Chemistry in Action is a comprehensive guide to the CMIS standard and related ECM concepts. In it, you'll find clear teaching and instantly useful examples for building content-centric client and server-side

applications that run against any CMIS-compliant repository. In fact, using the CMIS Workbench and the InMemory Repository from Apache Chemistry, you'll have running code talking to a real CMIS server by the end of chapter 1. This book requires some familiarity with content management systems and a standard programming language like Java or C#. No exposure to CMIS or Apache Chemistry is assumed. Purchase of the print book includes a free eBook in PDF, Kindle, and ePub formats from Manning Publications. What's Inside The only CMIS book endorsed by OASIS Complete coverage of the CMIS 1.0 and 1.1 specifications Cookbook-style tutorials and real-world examples About the Authors Florian Müller, Jay Brown, and Jeff Potts are among the original authors, contributors, and leaders of Apache Chemistry and the OASIS CMIS specification. They continue to shape CMIS implementations at Alfresco, IBM, and SAP.

Table of Contents PART 1 UNDERSTANDING CMIS  
Introducing CMIS Exploring the CMIS domain model  
Creating, updating, and deleting objects with CMIS  
CMIS metadata: types and properties Query PART 2 HANDS-ON  
CMIS CLIENT DEVELOPMENT Meet your new project: The Blend  
The Blend: read and query functionality The Blend: create, update, and delete functionality  
Using other client libraries Building mobile apps with CMIS PART 3  
ADVANCED TOPICS CMIS bindings Security and control Performance Building a CMIS server

This book is aimed at developers, designers, and architects who would like to build big data enterprise search solutions for their customers or organizations. No prior knowledge of Apache Hadoop and Apache Solr/Lucene technologies is required.

This book is part of Packt's Cookbook series; each chapter looks at a different aspect of working with Apache Solr. The recipes deal with common problems of working with Solr by

using easy-to-understand, real-life examples. The book is not in any way a complete Apache Solr reference and you should see it as a helping hand when things get rough on your journey with Apache Solr. Developers who are working with Apache Solr and would like to know how to combat common problems will find this book of great use. Knowledge of Apache Lucene would be a bonus but is not required. Big data has presented a number of opportunities across industries. With these opportunities come a number of challenges associated with handling, analyzing, and storing large data sets. One solution to this challenge is cloud computing, which supports a massive storage and computation facility in order to accommodate big data processing. *Managing and Processing Big Data in Cloud Computing* explores the challenges of supporting big data processing and cloud-based platforms as a proposed solution. Emphasizing a number of crucial topics such as data analytics, wireless networks, mobile clouds, and machine learning, this publication meets the research needs of data analysts, IT professionals, researchers, graduate students, and educators in the areas of data science, computer programming, and IT development.

[Copyright: 37857e0a08f125e37e908237f7ceacc7](https://www.dreambooks.com/book/37857e0a08f125e37e908237f7ceacc7)