# Apache Solr High Performance

Boost the performance of Solr instances and troubleshoot real-time problems

Surendra Mohan

# Apache Solr High Performance

Boost the performance of Solr instances and troubleshoot real-time problems

**Surendra Mohan**

[PACKT] PUBLISHING

open source*
community experience distilled

BIRMINGHAM - MUMBAI

# Apache Solr High Performance

# Credits

**Author**
Surendra Mohan

**Reviewers**
Azaz Desai
Ankit Jain
Mark Kerzner
Ruben Teijeiro

**Acquisition Editor**
Neha Nagwekar

**Content Development Editor**
Poonam Jain

**Technical Editor**
Krishnaveni Haridas

**Copy Editors**
Mradula Hegde
Alfida Paiva
Adithi Shetty

**Project Coordinator**
Puja Shukla

**Proofreaders**
Simran Bhogal
Ameesha Green
Maria Gould

**Indexers**
Monica Ajmera Mehta
Mariammal Chettiyar

**Graphics**
Abhinash Sahu

**Production Coordinator**
Saiprasad Kadam

**Cover Work**
Saiprasad Kadam

# About the Author

**Surendra Mohan**, who has served a few top-notch software organizations in varied roles, is currently a freelance software consultant. He has been working on various cutting-edge technologies such as Drupal and Moodle for more than nine years. He also delivers technical talks at various community events such as Drupal meet-ups and Drupal camps. To know more about him, his write-ups, and technical blogs, and much more, log on to `http://www.surendramohan.info/`.

He has also authored the book *Administrating Solr*, *Packt Publishing*, and has reviewed other technical books such as *Drupal 7 Multi Sites Configuration* and *Drupal Search Engine Optimization*, *Packt Publishing*, and titles on Drupal commerce and ElasticSearch, Drupal-related video tutorials, a title on Opsview, and many more.

# About the Reviewers

**Azaz Desai** has more than three years of experience in Mule ESB, jBPM, and Liferay technology. He is responsible for implementing, deploying, integrating, and optimizing services and business processes using ESB and BPM tools. He was a lead writer of *Mule ESB Cookbook*, *Packt Publishing*, and also played a vital role as a trainer on ESB. He currently provides training on Mule ESB to global clients. He has done various integrations of Mule ESB with Liferay, Alfresco, jBPM, and Drools. He was part of a key project on Mule ESB integration as a messaging system. He has worked on various web services and standards and frameworks such as CXF, AXIS, SOAP, and REST.

**Ankit Jain** holds a bachelor's degree in Computer Science Engineering from RGPV University, Bhopal, India. He has three years of experience in designing and architecting solutions for the Big Data domain and has been involved with several complex engagements. His technical strengths include Hadoop, Storm, S4, HBase, Hive, Sqoop, Flume, ElasticSearch, Machine Learning, Kafka, Spring, Java, and J2EE.

He also shares his thoughts on his personal blog at `http://ankitasblogger.blogspot.in/`. You can follow him on Twitter at `@mynameisanky`. He spends most of his time reading books and playing with different technologies. When not at work, Ankit spends time with his family and friends, watching movies, and playing games.

> I would like to thank my parents and brother for always being there for me.

**Mark Kerzner** holds degrees in Law, Maths, and Computer Science. He has been designing software for many years and Hadoop-based systems since 2008. He is the President of SHMsoft, a provider of Hadoop applications for various verticals, and a cofounder of the Hadoop Illuminated training and consulting, as well as the coauthor of the *Hadoop Illuminated* open source book. He has authored and coauthored several books and patents.

**Ruben Teijeiro** is an experienced frontend and backend web developer who had worked with several PHP frameworks for over a decade. His expertise is focused now on Drupal, with which he had collaborated in the development of several projects for some important organizations such as UNICEF and Telefonica in Spain and Ericsson in Sweden.

As an active member of the Drupal community, you can find him contributing to Drupal core, helping and mentoring other contributors, and speaking at Drupal events around the world. He also loves to share all that he has learned by writing in his blog, `http://drewpull.com`.

# www.PacktPub.com

## Support files, eBooks, discount offers and more

You might want to visit `www.PacktPub.com` for support files and downloads related to your book.

Did you know that Packt offers eBook versions of every book published, with PDF and ePub files available? You can upgrade to the eBook version at `www.PacktPub.com` and as a print book customer, you are entitled to a discount on the eBook copy. Get in touch with us at `service@packtpub.com` for more details.

At `www.PacktPub.com`, you can also read a collection of free technical articles, sign up for a range of free newsletters and receive exclusive discounts and offers on Packt books and eBooks.



`http://PacktLib.PacktPub.com`

Do you need instant solutions to your IT questions? PacktLib is Packt's online digital book library. Here, you can access, read and search across Packt's entire library of books.

### Why Subscribe?
- Fully searchable across every book published by Packt
- Copy and paste, print and bookmark content
- On demand and accessible via web browser

### Free Access for Packt account holders

If you have an account with Packt at `www.PacktPub.com`, you can use this to access PacktLib today and view nine entirely free books. Simply use your login credentials for immediate access.

# Table of Contents