# DATA QUALITY

The Kluwer International Series on
# ADVANCES IN DATABASE SYSTEMS

Series Editor
# Ahmed K. Elmagarmid

*Purdue University*
*West Lafayette, IN 47907*

***Other books in the Series:***

**THE FRACTAL STRUCTURE OF DATA REFERENCE:** *Applications to the Memory Hierarchy, Bruce McNutt;* ISBN: 0-7923-7945-4

**SEMANTIC MODELS FOR MULTIMEDIA DATABASE SEARCHING AND BROWSING,** *Shu-Ching Chen, R.L. Kashyap, and Arif Ghafoor;* ISBN: 0-7923-7888-1

**INFORMATION BROKERING ACROSS HETEROGENEOUS DIGITAL DATA: A Metadata-based Approach,** *Vipul Kashyap, Amit Sheth;* ISBN: 0-7923-7883-0

**DATA DISSEMINATION IN WIRELESS COMPUTING ENVIRONMENTS,** *Kian-Lee Tan and Beng Chin Ooi;* ISBN: 0-7923-7866-0

**MIDDLEWARE NETWORKS: Concept, Design and Deployment of Internet Infrastructure,** *Michah Lerner, George Vanecek, Nino Vidovic, Dad Vrsalovic;* ISBN: 0-7923-7840-7

**ADVANCED DATABASE INDEXING,** *Yannis Manolopoulos, Yannis Theodoridis, Vassilis J. Tsotras;* ISBN: 0-7923-7716-8

**MULTILEVEL SECURE TRANSACTION PROCESSING,** *Vijay Atluri, Sushil Jajodia, Binto George* ISBN: 0-7923-7702-8

**FUZZY LOGIC IN DATA MODELING,** *Guoqing Chen* ISBN: 0-7923-8253-6

**INTERCONNECTING HETEROGENEOUS INFORMATION SYSTEMS,** *Athman Bouguettaya, Boualem Benatallah, Ahmed Elmagarmid* ISBN: 0-7923-8216-1

**FOUNDATIONS OF KNOWLEDGE SYSTEMS: With Applications to Databases and Agents,** *Gerd Wagner* ISBN: 0-7923-8212-9

**DATABASE RECOVERY,** *Vijay Kumar, Sang H. Son* ISBN: 0-7923-8192-0

**PARALLEL, OBJECT-ORIENTED, AND ACTIVE KNOWLEDGE BASE SYSTEMS,** *Ioannis Vlahavas, Nick Bassiliades* ISBN: 0-7923-8117-3

**DATA MANAGEMENT FOR MOBILE COMPUTING,** *Evaggelia Pitoura, George Samaras* ISBN: 0-7923-8053-3

**MINING VERY LARGE DATABASES WITH PARALLEL PROCESSING,** *Alex A. Freitas, Simon H. Lavington* ISBN: 0-7923-8048-7

**INDEXING TECHNIQUES FOR ADVANCED DATABASE SYSTEMS,** *Elisa Bertino, Beng Chin Ooi, Ron Sacks-Davis, Kian-Lee Tan, Justin Zobel, Boris Shidlovsky, Barbara Catania* ISBN: 0-7923-9985-4

**INDEX DATA STRUCTURES IN OBJECT-ORIENTED DATABASES,** *Thomas A. Mueck, Martin L. Polaschek* ISBN: 0-7923-9971-4

# DATA QUALITY

*by*

**Richard Y. Wang**

**Mostapha Ziad**

**Yang W. Lee**

Visit Kluwer Online at:           http: www.kluweronline.com
and Kluwer's eBookstore at:       http: www.ebooks.kluweronline.com

*To our families ...*

# Table of Contents

*If you would not be forgotten,*
*As soon as you are dead and rotten,*
*Either write things worth reading,*
*Or do things worth the writing.*

Benjamin Franklin

This book provides an exposé of research and practice in data quality for technically oriented readers. It is based on the research conducted at the MIT Total Data Quality Management (TDQM) program and the work of other leading research institutions. It is intended for researchers, practitioners, educators and graduate students in the fields of Computer Science, Information Technology, and other inter-disciplinary fields. This book describes some of the pioneering research results that form a theoretical foundation for dealing with advanced issues related to data quality. In writing this book, our goal was to provide an overview of the cumulated research results from the MIT TDQM research perspective as it relates to database research. As such, it can be used by Ph.D. candidates who wish to further pursue their research in the data quality area and by IT professionals who wish to gain an insight into theoretical results and apply them in practice. This book also complements the authors' other well-received book *Quality Information and Knowledge* that deals with more managerially- and practice-oriented topics (Prentice Hall, 1999) and the book *Journey to Data Quality: A Roadmap for Higher Productivity* (in preparation).

The book is organized as follows. We first introduce fundamental concepts and a framework for total data quality management (TDQM) that encompasses data quality definition, measurement, analysis and improvement. We then present the database-oriented work from the TDQM program, work that is focused on the data quality area in the database field. Next, we profile research projects from leading research institutions. These projects broaden the reader's perspective and offer the reader a snapshot of some of the cutting-edge research projects. Finally, concluding remarks and future directions are presented.

*Quality.* Chapter 4 is based, in part, on the M.I.T. CISL working paper, *A Knowledge-Based Approach to Assisting In Data Quality Judgment,* co-authored by Yeona Jang, Henry Kon, and Richard Wang. Chapter 5 is based, in part, on the M.I.T. TDQM working paper, *A Data Quality Algebra for Estimating Query Result Quality.*

The MIT Context Interchange Project in Chapter 6 has been contributed by Stuart Madnick based on the paper, Metadata Jones and the Tower of Babel: The Challenge of Large-Scale Semantic Heterogeneity, published in the Proceedings of 1999 IEEE Meta-Data Conference. The European Union Data Warehouse Quality Project in Chapter 7 has been contributed by Matthias Jarke and Yannis Vassiliou based on their paper "Data Warehouse Quality: A Review of the DWQ Project", published in the Proceedings of the 1997 Conference on Information Quality. The Purdue University Data Quality Project in Chapter 8 has been written by Vassilios S. Verykios under the direction of Ahmed Elmagarmid. It profiles one aspect of the ongoing Purdue Univesity's data quality initiative. Originally, additional chapters based on the work conducted at Telcordia (formerly Bellcore), Georgia Mason University, The University of Queensland, Australia, the University of St. Gallen in Switzerland, Universidad de Buenos Aires in Argentina, and other institutions were planned. Unfortunately, we were unable to incorporate them into the current edition, but we certainly hope to do so in a future publication.

We appreciate Elizabeth Ziad for her love, support, and understanding, our children, Fori Wang, and Leyla, Ahmed, and Nasim Ziad who bring so much joy and happiness to our lives, as well as Abdallah, Abdelkrim, Ali and the Ziad family. In addition, we thank A. Russell and Arlene Lucid, Dris Djermoun, Kamel Youcef-Toumi, Hassan Raffa, Said Naili, Youcef Bennour, Boualem Kezim, Mohamed Gouali, Ahmed Sidi Yekhlef, Youcef Boudeffa, and Nacim Zeghlache for their help, support, and generosity. We also thank Jack and Evelyn Putman, and John, Mary Lou, Patrick, and Helen McCarthy for their love and support.

Last, but not least, we would like to thank our parents who instilled in us the love of learning.

*Richard Y. Wang*
Boston University, Boston, Massachusetts
& MIT TDQM program (http://web.mit.edu/tdqm/)
rwang@bu.edu

*Mostapha Ziad*
Suffolk University, Boston, Massachusetts
mziad@acad.suffolk.edu

*Yang W. Lee*
Northeastern University, Boston, Massachusetts
y.lee@nunet.neu.edu