# Advanced Information and Knowledge Processing

Jason T.L. Wang, Mohammed J. Zaki,
Hannu T.T. Toivonen and Dennis Shasha (Eds)

# Data Mining in Bioinformatics

With 110 Figures

Springer

Jason T.L. Wang, PhD
New Jersey Institute of Technology, USA

Mohammed J. Zaki, PhD
Computer Science Department, Rensselaer Polytechnic Institute, USA

Hannu T.T. Toivonen, PhD
University of Helsinki and Nokia Research Center

Dennis Shasha, PhD
New York University, USA

*Series Editors*
Xindong Wu
Lakhmi Jain

# Contents

# Contributors

**Peter Bajcsy**
Center for Supercomputing
    Applications
University of Illinois at
    Urbana-Champaign
USA

**Deb Bardhan**
Department of Computer Science
Rensselaer Polytechnic Institute
USA

**Chris Bystroff**
Department of Biology
Rensselaer Polytechnic Institute
USA

**Mukund Deshpande**
Oracle Corporation
USA

**Cinzia Di Pietro**
School of Medicine
University of Catania
Italy

**Alfredo Ferro**
Department of Mathematics and
    Computer Science
University of Catania
Italy

**Laurie Jane Hammel**
Department of Defense
USA

**Jiawei Han**
Department of Computer Science
University of Illinois at
    Urbana-Champaign
USA

**Kai Huang**
Department of Biological Sciences
Carnegie Mellon University
USA

**Donald P. Huddler**
Biophysics Research Division
University of Michigan
USA

**George Karypis**
Department of Computer Science
    and Engineering
University of Minnesota
USA

**Michihiro Kuramochi**
Department of Computer Science
    and Engineering
University of Minnesota
USA

**Lei Liu**
Center for Comparative
    and Functional Genomics
University of Illinois at
    Urbana-Champaign
USA

**Heikki Mannila**
Department of Computer Science
Helsinki University of Technology
Finland

**Robert F. Murphy**
Departments of Biological Sciences
    and Biomedical Engineering
Carnegie Mellon University
USA

**Vinay Nadimpally**
Department of Computer Science
Rensselaer Polytechnic Institute
USA

**Päivi Onkamo**
Department of Computer Science
University of Helsinki
Finland

**Roderic D. M. Page**
Division of Environmental
    and Evolutionary Biology
Institute of Biomedical and
    Life Sciences
University of Glasgow
United Kingdom

**Jignesh M. Patel**
Electrical Engineering and
    Computer Science Department
University of Michigan
USA

**Giuseppe Pigola**
Department of Mathematics and
    Computer Science
University of Catania
Italy

**Alfredo Pulvirenti**
Department of Mathematics and
    Computer Science
University of Catania
Italy

**Michele Purrello**
School of Medicine
University of Catania
Italy

**Marco Ragusa**
School of Medicine
University of Catania
Italy

**Marko Salmenkivi**
Department of Computer Science
University of Helsinki
Finland

**Petteri Sevon**
Department of Computer Science
University of Helsinki
Finland

**Dennis Shasha**
Courant Institute of Mathematical
    Sciences
New York University
USA

**Ambuj K. Singh**
Department of Computer Science
University of California at
    Santa Barbara
USA

**Hannu T. T. Toivonen**
Department of Computer Science
University of Helsinki
Finland

**Jason T. L. Wang**
Department of Computer Science
New Jersey Institute of Technology
USA

**Jiong Yang**
Department of Computer Science
University of Illinois at
    Urbana-Champaign
USA

**Mohammed J. Zaki**
Department of Computer Science
Rensselaer Polytechnic Institute
USA

**Kaizhong Zhang**
Department of Computer Science
University of Western Ontario
Canada

# Part I

# Overview