# Lecture Notes in Bioinformatics 4316

Subseries of Lecture Notes in Computer Science

Mehmet M. Dalkilic   Sun Kim
Jiong Yang (Eds.)

# Data Mining
# and Bioinformatics

First International Workshop, VDMB 2006
Seoul, Korea, September 11, 2006
Revised Selected Papers

Springer

# Preface

This volume contains the papers presented at the inaugural workshop on Data Mining and Bioinformatics at the 32nd International Conference on Very Large Data Bases (VLDB). The purpose of this workshop was to begin bringing together researchers from database, data mining, and bioinformatics areas to help leverage respective successes in each to the others. We also hope to expose the richness, complexity, and challenges in this area that involves mining very large complex biological data that will only grow in size and complexity as genome-scale high-throughput techniques become more routine. The problems are sufficiently different enough from traditional data mining problems (outside of life sciences) that novel approaches must be taken to data mine in this area. The workshop was held in Seoul, Korea, on September 11, 2006.

We received 30 submissions in response to the call for papers. Each submission was assigned to at least three members of the Program Committee. The Program Committee discussed the submission electronically, judging them on their importance, originality, clarity, relevance, and appropriateness to the expected audience. The Program Committee selected 15 papers for presentation. These papers are in the areas of microarray data analysis, bioinformatics system and text retrieval, application of gene expression data, and sequence analysis. Because of the format of the workshop and the high number of submissions, many good papers could not be included. Complementing the contributed papers, the program of VDMB 2006 included an invited talk by Simon Mercer, Program Manager for External Research, with an empahsis on life sciences.

We would like to thank the members of the Program Committee for their hard and expert work. We would also like to thank the VLDB organizers, the external reviewers, the authors, and the participants for their contribution to the continuing success of the workshop. Thanks also to Indiana University School of Informatics for the generous financial support.

October 2006

Mehmet Dalkilic  
Sun Kim  
Jiong Yang  
Program Chairs  
VDMB 2006

# VDMB 2006 Organization

## Program Committee Chairs

Mehmet Dalkilic (Indiana University, USA)
Sun Kim (Indiana University, USA)
Jiong Yang (Indiana University, USA)

## Program Committee Members

Mark Adams (Case Western University, USA)
Xue-wen Chen (University of Kansas, USA)
Jong Bhak (Korea Bioinformatics Center, Korea)
Dan Fay (Microsoft/Director Technical Computing, North America)
Hwan-Gue Cho (Busan National University, Korea)
Jeong-Hyeon Choi (Indiana University, USA)
Tony Hu (Drexel University, USA)
Jaewoo Kang (North Carolina State University, USA)
George Karypis (University of Minnesota, USA)
Doheon Lee (KAIST, Korea)
Jing Li (Case Western Reserve University, USA)
Yanda Li (Tsinghua University, China)
Birong Liao (Eli Lilly, USA)
Li Liao (University of Delaware, USA)
Huiqing Liu (University of Georgia, USA)
Lei Liu (University of Illinois at Urbana Champaign, USA)
Xiangjun (Frank) Liu (Tsinghua University, China)
Qingming Luo (Huazhong University, China)
Simon Mercer (Microsoft, USA)
Jian Pei (Simon Fraser University, Canada)
Meral Ozsoyoglu (Case Western Reserve University, USA)
Predrag Radivojac (Indiana University, USA)
Tetsuo Shibuya (University of Tokyo, Japan)
Keiji Takamoto (Case Western Reserve University, USA)
Haixu Tang (Indiana University, USA)
Anthony Tung (National University of Singapore, Singapore)
Wei Wang (University of North Carolina at Chapel Hill, USA)
Mohammed Zaki (Rensselaer Polytechnic Institute, USA)
Aidong Zhang (State University of New York at Buffalo, USA)

# Table of Contents