Roderic Guigó   Dan Gusfield (Eds.)

# Algorithms
# in Bioinformatics

Second International Workshop, WABI 2002
Rome, Italy, September 17-21, 2002
Proceedings

Springer

# Preface

We are pleased to present the proceedings of the *Second Workshop on Algorithms in Bioinformatics (WABI 2002)*, which took place on September 17-21, 2002 in Rome, Italy. The WABI workshop was part of a three-conference meeting, which, in addition to WABI, included the ESA and APPROX 2002. The three conferences are jointly called ALGO 2002, and were hosted by the Faculty of Engineering, University of Rome "La Sapienza". See `http://www.dis.uniroma1.it/~algo02` for more details.

The Workshop on Algorithms in Bioinformatics covers research in all areas of algorithmic work in bioinformatics and computational biology. The emphasis is on discrete algorithms that address important problems in molecular biology, genomics, and genetics, that are founded on sound models, that are computationally efficient, and that have been implemented and tested in simulations and on real datasets. The goal is to present recent research results, including significant work in progress, and to identify and explore directions of future research.

Original research papers (including significant work in progress) or state-of-the-art surveys were solicited on all aspects of algorithms in bioinformatics, including, but not limited to: exact and approximate algorithms for genomics, genetics, sequence analysis, gene and signal recognition, alignment, molecular evolution, phylogenetics, structure determination or prediction, gene expression and gene networks, proteomics, functional genomics, and drug design.

We received 83 submissions in response to our call for papers, and were able to accept about half of the submissions. In addition, WABI hosted two invited, distinguished lectures, given to the entire ALGO 2002 conference, by Dr. Ehud Shapiro of the Weizmann Institute and Dr. Gene Myers of Celera Genomics. An abstract of Dr. Shapiro's lecture, and a full paper detailing Dr. Myers lecture, are included in these proceedings.

We would like to sincerely thank all the authors of submitted papers, and the participants of the workshop. We also thank the program committee for their hard work in reviewing and selecting the papers for the workshop. We were fortunate to have on the program committee the following distinguished group of researchers:

Pankaj Agarwal (GlaxoSmithKline Pharmaceuticals, King of Prussia)
Alberto Apostolico (Università di Padova and Purdue University, Lafayette)
Craig Benham (University of California, Davis)
Jean-Michel Claverie (CNRS-AVENTIS, Marseille)
Nir Friedman (Hebrew University, Jerusalem)
Olivier Gascuel (Université de Montpellier II and CNRS, Montpellier)
Misha Gelfand (IntegratedGenomics, Moscow)
Raffaele Giancarlo (Università di Palermo)

David Gilbert (University of Glasgow)
Roderic Guigo (Institut Municipal d'Investigacions Mèdiques,
    Barcelona, co-chair)
Dan Gusfield (University of California, Davis, co-chair)
Jotun Hein (University of Oxford)
Inge Jonassen (Universitetet i Bergen)
Giuseppe Lancia (Università di Padova)
Bernard M.E. Moret (University of New Mexico, Albuquerque)
Gene Myers (Celera Genomics, Rockville)
Christos Ouzonis (European Bioinformatics Institute, Hinxton Hall)
Lior Pachter (University of California, Berkeley)
Knut Reinert (Celera Genomics, Rockville)
Marie-France Sagot (Université Claude Bernard, Lyon)
David Sankoff (Université de Montréal)
Steve Skiena (State University of New York, Stony Brook)
Gary Stormo (Washington University, St. Louis)
Jens Stoye (Universität Bielefeld)
Martin Tompa (University of Washington, Seattle)
Alfonso Valencia (Centro Nacional de Biotecnología, Madrid)
Martin Vingron (Max-Planck-Institut für Molekulare Genetik, Berlin)
Lusheng Wang (City University of Hong Kong)
Tandy Warnow (University of Texas, Austin)

We also would like to thank the WABI steering committee, Olivier Gascuel, Jotun Hein, Raffaele Giancarlo, Erik Meineche-Schmidt, and Bernard Moret, for inviting us to co-chair this program committee, and for their help in carrying out that task.

We are particularly indebted to Terri Knight of the University of California, Davis, Robert Castelo of the Universitat Pompeu Fabra, Barcelona, and Bernard Moret of the University of New Mexico, Albuquerque, for the extensive technical and advisory help they gave us. We could not have managed the reviewing process and the preparation of the proceedings without their help and advice.

Thanks again to everyone who helped to make WABI 2002 a success. We hope to see everyone again at WABI 2003.


July, 2002                                          Roderic Guigó and Dan Gusfield

# Table of Contents