# Lecture Notes in Bioinformatics 2812

Subseries of Lecture Notes in Computer Science

Gary Benson   Roderic Page (Eds.)

# Algorithms in Bioinformatics

Third International Workshop, WABI 2003
Budapest, Hungary, September 15-20, 2003
Proceedings

Springer

Series Editors

Sorin Istrail, Celera Genomics, Applied Biosystems, Rockville, MD, USA
Pavel Pevzner, University of California, San Diego, CA, USA
Michael Waterman, University of Southern California, Los Angeles, CA, USA

Volume Editors

Gary Benson
The Mount Sinai School of Medicine
Department of Biomathematical Sciences
Box 1023, One Gustave L. Levy Place, New York, NY 10029-6574, USA
E-mail: benson@camelot.mssm.edu

Roderic Page
University of Glasgow, Institute of Biomedical and Life Sciences
Division of Environmental and Evolutionary Biology
Glasgow G12 8QQ, Scotland
E-mail: r.page@bio.gla.ac.uk

# Preface

We are pleased to present the proceedings of the *Third Workshop on Algorithms in Bioinformatics (WABI 2003)*, which took place on September 15–20, 2003 in Budapest, Hungary. The WABI workshop was part of the four-conference meeting, ALGO 2003, which was locally organized by Dr. János Csirik, Head of the Department of Computer Science, József Attila University, Budapest. See `http://www.conferences.hu/ALGO2003/algo_2003.htm` for more details.

WABI focuses on discrete algorithms that address important problems in molecular biology, genomics, and genetics, that are founded on sound models, that are computationally efficient, that have been implemented and tested in simulations and on real datasets, and that provide new biological results. The workshop goals are to present recent research and identify and explore directions for future research.

We received 78 submissions in response to the call for papers and 36 were accepted. We would like to sincerely thank the authors of all submitted papers and the conference participants. We especially thank a terrific program committee for their diligent and thorough work in reviewing and selecting the papers. We were fortunate to have on the program committee the following distinguished group of researchers:

Amihood Amir (Bar Ilan University, Israel)
Alberto Apostolico (Purdue University)
Pierre Baldi (University of California, Irvine)
Gary Benson (Mount Sinai School of Medicine, New York; Co-chair)
Benny Chor (Tel Aviv University)
Nadia El-Mabrouk (University of Montreal)
Olivier Gascuel (LIRMM-CNRS, Montpellier)
Raffaele Giancarlo (Università di Palermo)
David Gilbert (University of Glasgow)
Jan Gorodkin (The Royal Veterinary and Agricultural University, Denmark)
Roderic Guigó (Institut Municipal d'Investigacions Mèdiques, Barcelona)
Dan Gusfield (University of California, Davis)
Jotun Hein (University of Oxford)
Daniel Huson (Tübingen University)
Simon Kasif (Boston University)
Gregory Kucherov (INRIA-Lorraine/LORIA)
Gad Landau (University of Haifa)
Thierry Lecroq (Université de Rouen)
Bernard M.E. Moret (University of New Mexico, Albuquerque)
Vincent Moulton (Uppsala University, Sweden)
Roderic Page (University of Glasgow; Co-chair)
Sophie Schbath (INRA, Jouy-en-Josas)

Charles Semple (University of Canterbury, New Zealand)
Jens Stoye (Universität Bielefeld)
Fengzhu Sun (University of Southern California)
Alfonso Valencia (Centro Nacional de Biotecnología, Madrid)
Jacques Van Helden (Université Libre de Bruxelles)
Louxin Zhang (National University of Singapore)

The program committee's work was greatly assisted by the helpful reviews provided by Federico Abascal, Ali Al-Shahib, Rumen Andonov, Ora Arbell, Sebastian Böcker, David Bryant, Peter Calabrese, Robert Castelo, Kwok Pui Choi, Miklós Csürös, Minghua Deng, Tobias Dezulian, Nadav Efraty, P.L. Erdos, Revital Eres, Eleazer Eskin, Jose Maria Fernandez, Pierre Flener, Jakob Fredslund, Dan Gieger, Robert Giegerich, Vladimir Grebinskiy, Pawel Herzyk, Mark Hoebeke, Robert Hoffmann, Katharina Huber, Michael Kaufmann, Carmel Kent, Jens Lagergren, Yinglei Lai, Xiaoman Li, Gerton Lunter, Rune Lyngsoe, Laurent Mouchard, Pierre Nicolas, Laurent Noe, Sebastian Oehm, Christian N.S. Pedersen, Johann Pelfrene, Shalom Rackovsky, Mathieu Raffinot, Kim Roland Rasmussen, Christian Rausch, Knut Reinert, Olivier Sand, Klaus-Bernd Schürmann, Steven Skiena, Dina Sokol, Y.S. Song, W. Szpankowski, Helene Touzet, Michael Tress, Juris Viksna, Lusheng Wang, Zohar Yakhini, Kui Zhang, and Michal Ziv-Ukelson.

We also thank the WABI steering committee, Olivier Gascuel, Raffaele Giancarlo, Roderic Guigó, Dan Gusfield, and Bernard Moret, for inviting us to co-chair the conference and for their help in carrying out that task.

We are particularly indebted to Kevin Kelliher of the Mount Sinai School of Medicine, New York, for his conscientious administration of the CyberChair software used to manage the review process, and Robert Castelo at Universitat Pompeu Fabra, Barcelona, for generously sharing the scripts he developed to generate the final copy of last year's WABI. Production of the proceedings was greatly assisted by Richard Koch's wonderful TeXShop software, and large quantities of Peet's coffee.

Thanks again to all who participated to make WABI 2003 a success. It has been, for us, a challenging and rewarding experience.

July 2003                                        Gary Benson and Roderic Page

# Table of Contents

# Pattern and Motif Discovery

# Phylogenetic Analysis

## Polymorphism

## Protein Structure

## Sequence Alignment

## String Algorithms