Olivier Gascuel  Bernard M.E. Moret (Eds.)

# Algorithms
# in Bioinformatics

First International Workshop, WABI 2001
Århus Denmark, August 28-31, 2001
Proceedings

Springer

# Preface

We are very pleased to present the proceedings of the *First Workshop on Bioinformatics (WABI 2001)*, which took place in Aarhus on August 28–31, 2001, under the auspices of the *European Association for Theoretical Computer Science (EATCS)* and the Danish Center for *Basic Research in Computer Science (BRICS)*.

The *Workshop on Algorithms in Bioinformatics* covers research on all aspects of algorithmic work in bioinformatics. The emphasis is on discrete algorithms that address important problems in molecular biology. These are founded on sound models, are computationally efficient, and 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. Specific topics of interest include, but are not limited to:

  – Exact and approximate algorithms for genomics, sequence analysis, gene and signal recognition, alignment, molecular evolution, structure determination or prediction, gene expression and gene networks, proteomics, functional genomics, and drug design.
  – Methods, software and dataset repositories for development and testing of such algorithms and their underlying models.
  – High-performance approaches to computationally hard problems in bioinformatics, particularly optimization problems.

A major goal of the workshop is to bring together researchers spanning the range from abstract algorithm design to biological dataset analysis, to encourage dialogue between application specialists and algorithm designers, mediated by algorithm engineers and high-performance computing specialists. We believe that such a dialogue is necessary for the progress of computational biology, inasmuch as application specialists cannot analyze their datasets without fast and robust algorithms and, conversely, algorithm designers cannot produce useful algorithms without being aware of the problems faced by biologists. Part of this mix was achieved automatically this year by colocating into a single large conference, *ALGO 2001*, three workshops: *WABI 2001*, the *5th Workshop on Algorithm Engineering (WAE 2001)*, and the *9th European Symposium on Algorithms (ESA 2001)*, and sharing keynote addresses among the three workshops. *ESA* attracts algorithm designers, mostly with a theoretical leaning, while *WAE* is explicitly targeted at algorithm engineers and algorithm experimentalists.

These proceedings reflect such a mix. We received over 50 submissions in response to our call and were able to accept 23 of them, ranging from mathematical tools through to experimental studies of approximation algorithms and reports on significant computational analyses. Numerous biological problems are dealt with, including genetic mapping, sequence alignment and sequence analysis, phylogeny, comparative genomics, and protein structure.

We were also fortunate to attract Dr. Gene Myers, Vice-President for Informatics Research at Celera Genomics, and Prof. Jotun Hein, Aarhus University, to address the joint workshops, joining five other distinguished speakers (Profs. Herbert Edelsbrunner and Lars Arge from Duke University, Prof. Susanne Albers from Dortmund University, Prof. Uri Zwick from Tel Aviv University, and Dr. Andrei Broder from Alta Vista). The quality of the submissions and the interest expressed in the workshop is promising – plans for next year's workshop are under way.

We would like to thank all the authors for submitting their work to the workshop and all the presenters and attendees for their participation. We were particularly fortunate in enlisting the help of a very distinguished panel of researchers for our program committee, which undoubtedly accounts for the large number of submissions and the high quality of the presentations. Our heartfelt thanks go to all:

Craig Benham (Mt Sinai School of Medicine, New York, USA)
Mikhail Gelfand (Integrated Genomics, Moscow, Russia)
Raffaele Giancarlo (U. di Palermo, Italy)
Michael Hallett (McGill U., Canada)
Jotun Hein (Aarhus U., Denmark)
Michael Hendy (Massey U., New Zealand)
Inge Jonassen (Bergen U., Norway)
Junhyong Kim (Yale U., New Haven, USA)
Jens Lagergren (KTH Stockholm, Sweden)
Edward Marcotte (U. Texas Austin, USA)
Satoru Miyano (Tokyo U., Japan)
Gene Myers (Celera Genomics, USA)
Marie-France Sagot (Institut Pasteur, France)
David Sankoff (U. Montreal, Canada)
Thomas Schiex (INRA Toulouse, France)
Joao Setubal (U. Campinas, Sao Paolo, Brazil)
Ron Shamir (Tel Aviv U., Israel)
Lisa Vawter (GlaxoSmithKline, USA)
Martin Vingron (Max Planck Inst. Berlin, Germany)
Tandy Warnow (U. Texas Austin, USA)

In addition, the opinion of several other researchers was solicited. These subreferees include Tim Beissbarth, Vincent Berry, Benny Chor, Eivind Coward, Ingvar Eidhammer, Thomas Faraut, Nicolas Galtier, Michel Goulard, Jacques van Helden, Anja von Heydebreck, Ina Koch, Chaim Linhart, Hannes Luz, Vsevolod Yu, Michal Ozery, Itsik Pe'er, Sven Rahmann, Katja Rateitschak, Eric Rivals, Mikhail A. Roytberg, Roded Sharan, Jens Stoye, Dekel Tsur, and Jian Zhang. We thank them all.

Lastly, we thank Prof. Erik Meineche-Schmidt, BRICS codirector, who started the entire enterprise by calling on one of us (Bernard Moret) to set up the workshop and who led the team of committee chairs and organizers through the

setup, development, and actual events of the three combined workshops, with the assistance of Prof. Gerth Brødal.

We hope that you will consider contributing to *WABI 2002*, through a submission or by participating in the workshop.

June 2001                              Olivier Gascuel and Bernard M.E. Moret

# Table of Contents