

Lecture Notes in Bioinformatics

2983

Edited by S. Istrail, P. Pevzner, and M. Waterman

Editorial Board: A. Apostolico S. Brunak M. Gelfand
T. Lengauer S. Miyano G. Myers M.-F. Sagot D. Sankoff
R. Shamir T. Speed M. Vingron W. Wong

Subseries of Lecture Notes in Computer Science

Springer

Berlin

Heidelberg

New York

Hong Kong

London

Milan

Paris

Tokyo

Sorin Istrail Michael Waterman
Andrew Clark (Eds.)

Computational Methods for SNPs and Haplotype Inference

DIMACS/RECOMB Satellite Workshop
Piscataway, NJ, USA, November 21-22, 2002
Revised Papers



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

Sorin Istrail
Celera Genomics, Applied Biosystems
45 West Gude Drive, Rockville, MD 20850, USA
E-mail: sorin.istrail@celera.com

Michael Waterman
Celera Genomics and
University of Southern California
1042W 36th Place, DRB 155, Los Angeles, CA 90089-1113, USA
E-mail: msw@hto.usc.edu

Andrew Clark
Celera Genomics and
Cornell University, Molecular Biology and Genetics
Ithaca, NY 14853, USA
E-mail: Andy.Clark@celera.com

Cataloging-in-Publication Data applied for

A catalog record for this book is available from the Library of Congress.

Bibliographic information published by Die Deutsche Bibliothek
Die Deutsche Bibliothek lists this publication in the Deutsche Nationalbibliografie;
detailed bibliographic data is available in the Internet at <<http://dnb.ddb.de>>.

CR Subject Classification (1998): H.2.8, F.2.2, G.2.1, G.3, H.3.3, I.5, J.3

ISSN 0302-9743
ISBN 3-540-21249-3 Springer-Verlag Berlin Heidelberg New York

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, re-use of illustrations, recitation, broadcasting, reproduction on microfilms or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer-Verlag. Violations are liable for prosecution under the German Copyright Law.

Springer-Verlag is a part of Springer Science+Business Media
springeronline.com

© Springer-Verlag Berlin Heidelberg 2004
Printed in Germany

Typesetting: Camera-ready by author, data conversion by Boller Mediendesign
Printed on acid-free paper SPIN: 10993057 06/3142 5 4 3 2 1 0

Preface

The 1st Computational Methods for SNPs and Haplotype Inference Workshop was held on November 21–22, 2002 at the DIMACS Center for Discrete Mathematics and Theoretical Computer Science.

The workshop focused on methods for SNP and haplotype analysis and their applications to disease associations. The ability to score large numbers of DNA variants (SNPs) in large samples of humans is rapidly accelerating, as is the demand to apply these data to tests of association with diseased states. The problem suffers from excessive dimensionality, so any means of reducing the number of dimensions of the space of genotype classes in a biologically meaningful way would likely be of benefit. Linked SNPs are often statistically associated with one another (in "linkage disequilibrium"), and the number of distinct configurations of multiple tightly linked SNPs in a sample is often far lower than one would expect from independent sampling. These joint configurations, or haplotypes, might be a more biologically meaningful unit, since they represent sets of SNPs that co-occur in a population. Recently there has been much excitement over the idea that such haplotypes occur as blocks across the genome, as these blocks suggest that fewer distinct SNPs need to be scored to capture the information about genotype identity. There is need for formal analysis of this dimension reduction problem, for formal treatment of the hierarchical structure of haplotypes, and for consideration of the utility of these approaches toward meeting the end goal of finding genetic variants associated with complex diseases.

The workshop featured the following invited speakers:

Peter Donnelly (Oxford University), Kathryn Roeder (Carnegie Mellon University), Jonathan Pritchard (University of Chicago), Molly Przeworski (Max Planck Institute), Maoxia Zheng (University of Chicago), Elizabeth Thompson (University of Washington), Monty Slatkin (University of California, Berkeley), Dahlia Nielsen (North Carolina State University), Matthew Stephens (University of Washington), Andrew Clark (Cornell University and Celera/Applied Biosystems), Sorin Istrail (Celera/Applied Biosystems), David Cutler (Johns Hopkins University), Magnus Nordborg (University of Southern California), Bruce Rannala (University of Alberta), Russell Schwartz (Carnegie Mellon University), Fengzhu Sun (University of Southern California), Jun Liu (Stanford University), Jinghui Zhang (National Cancer Institute, NIH), Dan Gusfield (University of California, Davis), Eran Halperin (University of California, Berkeley), Vineet Bafna (The Center for the Advancement of Genomics), David Altschuler (Harvard Medical School), Nancy Cox (University of Chicago), Francisco de la Vega (Applied Biosystems), Li Jin (University of Cincinnati) and Steve Sherry (National Center for Biotechnology Information, NIH).

This volume includes papers on which the presentations of Andrew Clark, Dan Gusfield, Sorin Istrail, Tianhua Niu, Jonathan Pritchard, Russell Schwartz, Elizabeth Thompson, Bruce Rannala, Fengzhu Sun, and Maoxia Zheng were based. It also includes the collection of abstracts of all the presentations.

We would like to thank Merissa Henry for outstanding workshop organization and editorial support. We would also like to thank the DIMACS Center for Discrete Mathematics and Theoretical Computer Science for providing financial support and excellent organization of the workshop.

January 2004

Andrew G. Clark
Sorin Istrail
Michael Waterman

Workshop Organizers

Table of Contents

Trisomic Phase Inference	1
<i>A.G. Clark, E.T. Dermitzakis, S.E. Antonarakis</i>	
An Overview of Combinatorial Methods for Haplotype Inference	9
<i>D. Gusfield</i>	
A Survey of Computational Methods for Determining Haplotypes	26
<i>B.V. Halldórsson, V. Bafna, N. Edwards, R. Lippert, S. Yooseph, S. Istrail</i>	
Haplotype Inference and Its Application in Linkage Disequilibrium Mapping	48
<i>T. Niu, X. Lu, H. Kang, Z.S. Qin, J.S. Liu</i>	
Inferring Piecewise Ancestral History from Haploid Sequences	62
<i>R. Schwartz, A.G. Clark, S. Istrail</i>	
Haplotype Blocks in Small Populations	74
<i>E.A. Thompson, N.H. Chapman</i>	
Simulating a Coalescent Process with Recombination and Ascertainment .	84
<i>Y. Wang, B. Rannala</i>	
Dynamic Programming Algorithms for Haplotype Block Partitioning and Tag SNP Selection Using Haplotype Data or Genotype Data	96
<i>K. Zhang, T. Chen, M.S. Waterman, Z.S. Qin, J.S. Liu, F. Sun</i>	
Parametric Bootstrap for Assessment of Goodness of Fit of Models for Block Haplotype Structure	113
<i>M. Zheng, M.S. McPeek</i>	
A Coalescent-Based Approach for Complex Disease Mapping	124
<i>S. Zöllner, J.K. Pritchard</i>	
Abstracts	
Haplotyping as Perfect Phylogeny	131
<i>V. Bafna, D. Gusfield, G. Lancia, S. Yooseph</i>	
Exhaustive Enumeration and Bayesian Phase Inference	132
<i>A. Clark</i>	
How Does Choice of Polymorphism Influence Estimation of LD and Mapping?	133
<i>N. Cox</i>	

VIII Table of Contents

Haplotype Inference in Random Population Samples	134
<i>D. Cutler</i>	
Bayesian Methods for Statistical Reconstruction of Haplotypes	135
<i>P. Donnelly</i>	
Combinatorial Approaches to Haplotype Inference	136
<i>D. Gusfield</i>	
Large Scale Recovery of Haplotypes from Genotype Data Using Imperfect Phylogeny	137
<i>E. Halperin</i>	
Haplotype Inference and Haplotype Information	138
<i>J. Liu</i>	
Multi-locus Linkage Disequilibrium and Haplotype-Based Tests of Association	139
<i>D. Nielsen</i>	
The Pattern of Polymorphism on Human Chromosome 21	140
<i>M. Nordborg</i>	
Use of a Local Approximation to the Ancestral Recombination Graph for Fine Mapping Disease Genes	141
<i>J. Pritchard, S. Zöllner</i>	
Insights into Recombination from Patterns of Linkage Disequilibrium	142
<i>M. Przeworski</i>	
Joint Bayesian Estimation of Mutation Location and Age Using Linkage Disequilibrium	143
<i>B. Rannala</i>	
Evolutionary-Based Association Analysis Using Haplotype Data	144
<i>K. Roeder</i>	
Inferring Piecewise Ancestral History from Haploid Sequences	145
<i>R. Schwartz</i>	
Testing for Differences in Haplotype Frequencies in Case-Control Studies .	146
<i>M. Slatkin</i>	
Haplotypes, Hotspots, and a Multilocus Model for Linkage Disequilibrium	147
<i>M. Stephens</i>	
Dynamic Programming Algorithms for Haplotype Block Partition and Applications to Association Studies	148
<i>F. Sun</i>	

Genome Sharing in Small Populations	149
<i>E. Thompson</i>	
Patterns of Linkage Disequilibrium across Human Chromosomes 6, 21, AND 22	150
<i>F. de la Vega</i>	
A Software System for Automated and Visual Analysis of Functionally Annotated Haplotypes	151
<i>J. Zhang</i>	
Assessment of Goodness of Fit of Models for Block Haplotype Structure .	152
<i>M. Zheng, M.S. McPeek</i>	
Author Index	153