Bin Ma
Kaizhong Zhang (Eds.)

# Combinatorial Pattern Matching

**18th Annual Symposium, CPM 2007**
**London, Canada, July 2007**
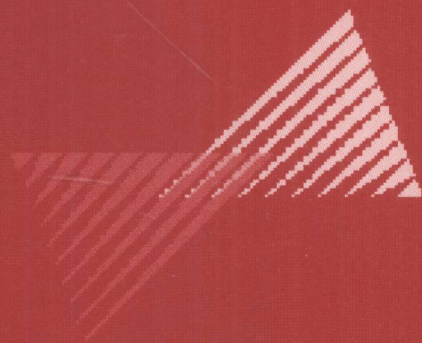**Proceedings**

Springer

Bin Ma   Kaizhong Zhang (Eds.)

# Combinatorial
# Pattern Matching

18th Annual Symposium, CPM 2007
London, Canada, July 9-11, 2007
Proceedings

◇ Springer

Volume Editors

Bin Ma
Kaizhong Zhang
University of Western Ontario
Department of Computer Science
London, Ontario, N6A 5B7, Canada
E-mail: {bma; kzhang}@csd.uwo.ca

# Lecture Notes in Computer Science 4580

# Preface

The papers contained in this volume were presented at the 18th Annual Symposium on Combinatorial Pattern Matching (CPM 2007) held at the University of Western Ontario, in London, Ontario, Canada from July 9 to 11, 2007.

All the papers presented at the conference are original research contributions on computational pattern matching and analysis, data compression and compressed text processing, suffix arrays and trees, and computational biology. They were selected from 64 submissions. Each submission was reviewed by at least three reviewers. The committee decided to accept 32 papers. The programme also included three invited talks by Tao Jiang from the University of California, Riverside, USA, S. Muthukrishnan from Rutgers University, USA, and Frances Yao from City University of Hong Kong, Hong Kong.

Combinatorial Pattern Matching addresses issues of searching and matching strings and more complicated patterns such as trees, regular expressions, graphs, point sets, and arrays. The goal is to derive non-trivial combinatorial properties of such structures and to exploit these properties in order to either achieve superior performance for the corresponding computational problems or pinpoint conditions under which searches cannot be performed efficiently.

The Annual Symposium on Combinatorial Pattern Matching started in 1990, and has since taken place every year. The objective of the annual CPM meetings is to provide an international forum for research in combinatorial pattern matching and related applications. Previous CPM meetings were held in Paris, London, Tucson, Padova, Asilomar, Helsinki, Laguna Beach, Aarhus, Piscataway, Warwick, Montreal, Jerusalem, Fukuoka, Morelia, Istanbul, Jeju Island, and Barcelona. Selected papers from the first meeting appeared in volume 92 of *Theoretical Computer Science*, from the 11th meeting in volume 2 of the *Journal of Discrete Algorithms*, from the 12th meeting in volume 146 of *Discrete Applied Mathematics*, and from the 14th meeting in volume 3 of the *Journal of Discrete Algorithms*. Starting from the 3rd meeting, the proceedings of all meetings have been published in the LNCS series, volumes 644, 684, 807, 937, 1075, 1264, 1448, 1645, 1848, 2089, 2373, 2676, 3109, 3537, 4009, and 4580.

The whole submission and review process, as well as the production of this volume, was carried out with the help of the EasyChiar system. The conference was sponsored by the University of Western Ontario and by the Fields Institute.

April 2007

Bin Ma
Kaizhong Zhang

# Conference Organization

## Programme Chairs

Bin Ma
Kaizhong Zhang

## Programme Committee

Tatsuya Akutsu
Amihood Amir
Raffaele Giancarlo
Gregory Kucherov
Ming Li
Guohui Lin
Heikki Mannila
Gonzalo Navarro
Ron Pinter
Mathieu Raffinot
Cenk Sahinalp
David Sankoff
Steve Skiena
James Storer
Masayuki Takeda
Gabriel Valiente
Martin Vingron
Lusheng Wang

## Local Organization

Meg Borthwick
Jingping Liu
Cheryl McGrath
Roberto Solis-Oba (Co-chair)
Kaizhong Zhang (Co-chair)

## External Reviewers

Cagri Aksay
Hideo Bannai
Petra Berenbrink

Guillaume Blin
Marie-Pierre Béal
Cedric Chauve

Shihyen Chen
Phuong Dao
Guillaume Fertin
Tom Friedetzky
Liliana Félix
Aris Gionis
Stefan Haas
Niina Haiminen
Iman Hajiresouliha
Tzvika Hartman
Rodrigo Hausen
Hannes Heikinheimo
Yasuto Higa
Fereydoun Hormozdiari
Lucian Ilie
Shunsuke Inenaga
Hossein Jowhari
Oren Kapah
Emre Karakoc
Orgad Keller
Takuya Kida
Roman Kolpakov
Tsvi Kopelowitz
Dennis Kostka
Juha Kärkkäinen
Gad Landau
Michael Lappe
Thierry Lecroq
Avivit Levy
Weiming Li

Yury Lifshits
Jingping Liu
Mercè Llabrés
Antoni Lozano
Giovanni Manzini
Conrado Martínez
Igor Nor
Pasi Rastas
Tobias Rausch
David Reese
Hugues Richard
Jairo Rocha
Oleg Rokhlenko
Francesc Rosselló
Dominique Rossin
Wojciech Rytter
Kunihiko Sadakane
Hiroshi Sakamoto
Rahaleh Salari
Marcel Schulz
Jouni Seppanen
Dina Sokol
Jens Stoye
Takeyuki Tamura
Eric Tannier
Helene Touzet
Antti Ukkonen
Tomas Vinar
Robert Warren
Lei Xin

# Lecture Notes in Computer Science

For information about Vols. 1–4492

please contact your bookseller or Springer

# Table of Contents

## Session 4: Computational Biology I

## Session 5: Computational Biology II

## Session 6: Algorithmic Techniques II

## Session 7: Data Compression II

## Session 8: Computational Biology III

## Session 9: Pattern Analysis

## Session 10: Suffix Arrays and Trees

# A Combinatorial Approach to Genome-Wide Ortholog Assignment: Beyond Sequence Similarity Search

Tao Jiang

Computer Science Department, University of California - Riverside
jiang@cs.ucr.edu

**Abstract.** The assignment of orthologous genes between a pair of genomes is a fundamental and challenging problem in comparative genomics. Existing methods that assign orthologs based on the similarity between DNA or protein sequences may make erroneous assignments when sequence similarity does not clearly delineate the evolutionary relationship among genes of the same families. In this paper, we present a new approach to ortholog assignment that takes into account both sequence similarity and evolutionary events at genome level, where orthologous genes are assumed to correspond to each other in the most parsimonious evolving scenario under genome rearrangement and gene duplication. It is then formulated as a problem of computing the signed reversal distance with duplicates between two genomes of interest, for which an efficient heuristic algorithm was constructed based on solutions to two new optimization problems, minimum common partition and maximum cycle decomposition. Following this approach, we have implemented a high-throughput system for assigning orthologs on a genome scale, called MSOAR, and tested it on both simulated data and real genome sequence data. Our predicted orthologs between the human and mouse genomes are strongly supported by ortholog and protein function information in authoritative databases, and predictions made by other key ortholog assignment methods such as Ensembl, Homologene, INPARANOID, and HGNC. The simulation results demonstrate that MSOAR in general performs better than the iterated exemplar algorithm of D. Sankoff's in terms of identifying true exemplar genes.

This is joint work with X. Chen (Nanyang Tech. Univ., Singapore), Z. Fu (UCR), J. Zheng (NCBI), V. Vacic (UCR), P. Nan (SCBIT), Y. Zhong (SCBIT), and S. Lonardi (UCR).

# Stringology: Some Classic and Some Modern Problems

S. Muthukrishnan

Department of Computer Science, Rutgers University
and
Google Inc.
muthu@cs.rutgers.edu

**Abstract.** We examine some of the classic problems related to suffix trees from 70's and show some recent results on sorting suffixes with small space and suffix selection. Further, we introduce modern versions of suffix sorting and their application to XML processing. The study of combinatorial aspects of strings continues to flourish, and we present several open problems with modern applications.

# Algorithmic Problems in Scheduling Jobs on Variable-Speed Processors

Frances F. Yao

Department of Computer Science,
City University of Hong Kong
Hong Kong SAR, China
csfyao@cityu.edu.hk

**Abstract.** Power and heat have become two of the major concerns for the computer industry, which struggles to cope with the energy and cooling costs for servers, as well as the short battery life of portable devices. Dynamic Voltage Scaling (DVS) has emerged as a useful technique: e.g. Intel's newest Foxton technology enables a chip to run at 64 different speed levels. Equipped with DVS technology, the operating system can then save CPU's energy consumption by scheduling tasks wisely. A schedule that finishes the given tasks within their timing constraints while using minimum total energy (among all feasible schedules) is called an optimal DVS schedule. A theoretical model for DVS scheduling was proposed in a paper by Yao, Demers and Shenker in 1995, along with a well-formed characterization of the optimum and an algorithm for computing it. This algorithm has remained as the most efficient known despite many investigations of this model. In this talk, we will first give an overview of the DVS scheduling problem, and then present the latest improved results for computing the optimal schedule in both the finite and the continuous (infinite speed levels) models. Related results on efficient on-line scheduling heuristics will also be discussed.

# Speeding Up HMM Decoding and Training by Exploiting Sequence Repetitions

Shay Mozes[1,*], Oren Weimann[1], and Michal Ziv-Ukelson[2,**]

[1] MIT Computer Science and Artificial Intelligence Laboratory,
32 Vassar Street, Cambridge, MA 02139, USA
`shaymozes@gmail.com,oweimann@mit.edu`
[2] School of Computer Science, Tel-Aviv University, Tel-Aviv 69978, Israel
`michaluz@post.tau.ac.il`

**Abstract.** We present a method to speed up the dynamic program algorithms used for solving the HMM decoding and training problems for discrete time-independent HMMs. We discuss the application of our method to Viterbi's decoding and training algorithms [21], as well as to the forward-backward and Baum-Welch [4] algorithms. Our approach is based on identifying repeated substrings in the observed input sequence. We describe three algorithms based alternatively on byte pair encoding (BPE) [19], run length encoding (RLE) and Lempel-Ziv (LZ78) parsing [22]. Compared to Viterbi's algorithm, we achieve a speedup of $\Omega(r)$ using BPE, a speedup of $\Omega(\frac{r}{\log r})$ using RLE, and a speedup of $\Omega(\frac{\log n}{k})$ using LZ78, where $k$ is the number of hidden states, $n$ is the length of the observed sequence and $r$ is its compression ratio (under each compression scheme). Our experimental results demonstrate that our new algorithms are indeed faster in practice. Furthermore, unlike Viterbi's algorithm, our algorithms are highly parallelizable.

**Keywords:** HMM, Viterbi, dynamic programming, compression.

## 1   Introduction

Over the last few decades, Hidden Markov Models (HMMs) proved to be an extremely useful framework for modeling processes in diverse areas such as error-correction in communication links [21], speech recognition [6], optical character recognition [2], computational linguistics [17], and bioinformatics [12].

The core HMM-based applications fall in the domain of classification methods and are technically divided into two stages: a training stage and a decoding stage. During the *training* stage, the emission and transition probabilities of an HMM are estimated, based on an input set of observed sequences. This stage is usually executed once as a preprocessing stage and the generated ("trained") models are stored in a database. Then, a *decoding* stage is run, again and again, in order to

---