

Alexander S. Kulikov
Sergei O. Kuznetsov
Pavel Pevzner (Eds.)

CERIST
LNCS 8486
BIBLIOTHEQUE DU

Combinatorial Pattern Matching

25th Annual Symposium, CPM 2014
Moscow, Russia, June 16–18, 2014
Proceedings



Commenced Publication in 1973

Founding and Former Series Editors:

Gerhard Goos, Juris Hartmanis, and Jan van Leeuwen

Editorial Board

David Hutchison

Lancaster University, UK

Takeo Kanade

Carnegie Mellon University, Pittsburgh, PA, USA

Josef Kittler

University of Surrey, Guildford, UK

Jon M. Kleinberg

Cornell University, Ithaca, NY, USA

Alfred Kobsa

University of California, Irvine, CA, USA

Friedemann Mattern

ETH Zurich, Switzerland

John C. Mitchell

Stanford University, CA, USA

Moni Naor

Weizmann Institute of Science, Rehovot, Israel

Oscar Nierstrasz

University of Bern, Switzerland

C. Pandu Rangan

Indian Institute of Technology, Madras, India

Bernhard Steffen

TU Dortmund University, Germany

Demetri Terzopoulos

University of California, Los Angeles, CA, USA

Doug Tygar

University of California, Berkeley, CA, USA

Gerhard Weikum

Max Planck Institute for Informatics, Saarbruecken, Germany

Alexander S. Kulikov Sergei O. Kuznetsov
Pavel Pevzner (Eds.)

Combinatorial Pattern Matching

25th Annual Symposium, CPM 2014
Moscow, Russia, June 16-18, 2014
Proceedings



Springer

Volume Editors

Alexander S. Kulikov
St. Petersburg Department of Steklov Institute of Mathematics
27 Fontanka
St. Petersburg 191023, Russia
E-mail: kulikov@logic.pdmi.ras.ru

Sergei O. Kuznetsov
National Research University Higher School of Economics
3 Bolshoy Trekhsvyatitskiy pereulok
Moscow 109028, Russia
E-mail: skuznetsov@hse.ru

Pavel Pevzner
University of California at San Diego
9500 Gilman Drive, EBU3b 4236
La Jolla, CA 92093-0404, USA
E-mail: ppevzner@cs.ucsd.edu

ISSN 0302-9743 e-ISSN 1611-3349
ISBN 978-3-319-07565-5 e-ISBN 978-3-319-07566-2
DOI 10.1007/978-3-319-07566-2
Springer Cham Heidelberg New York Dordrecht London

Library of Congress Control Number: 2014939431

LNCS Sublibrary: SL 1 – Theoretical Computer Science and General Issues

© Springer International Publishing Switzerland 2014

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Typesetting: Camera-ready by author, data conversion by Scientific Publishing Services, Chennai, India

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

Preface

This volume contains the papers presented at the 25th Annual Symposium on Combinatorial Pattern Matching held during June 15–17, 2014 in Moscow, Russia. The hosting university was National Research University Higher School of Economics.

There were 54 submissions from 25 countries. Each submission was reviewed by at least three Program Committee members. The committee decided to accept 28 papers, corresponding to an acceptance rate of 54%. We thank the members of the Program Committee and all additional external reviewers for their hard work that resulted in this excellent program. Their names are listed on the following pages. The whole submission and review process was carried out with the invaluable help of the EasyChair conference system.

The year 2014 marks the quarter-of-a-century milestone for the CPM symposium series. CPM 2014 celebrated this event with invited talks by the co-founders of the conference series: Alberto Apostolico (Georgia Institute of Technology, USA, and IASI-CNR, Italy), Maxime Crochemore (King's College London, UK, and Université Paris-Est, France), Zvi Galil (Georgia Institute of Technology, USA), and Udi Manber (Google, USA). In addition, CPM 2014 featured a special keynote lecture by Gene Myers (Max Planck Institute, Germany) that highlighted the contribution of the CPM community to bioinformatics.

The objective of the annual CPM meetings is to provide an international forum for research in combinatorial pattern matching and related applications. It 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 combinatorial properties of such structures and to exploit these properties in order to achieve superior performance for the corresponding computational problems. The meeting also deals with problems in computational biology, data compression and data mining, coding, information retrieval, natural language processing, and pattern recognition.

The Annual Symposium on Combinatorial Pattern Matching started in 1990, and has since taken place every year. Previous CPM meetings were held in Paris, London (UK), Tucson, Padova, Asilomar, Helsinki, Laguna Beach, Aarhus, Piscataway, Warwick, Montreal, Jerusalem, Fukuoka, Morelia, Istanbul, Jeju Island, Barcelona, London (Ontario, Canada), Pisa, Lille, New York, Palermo, Helsinki, Bad Herrenalb. This year's meeting was the first in Russia. Starting from the third meeting, proceedings of all meetings have been published in the LNCS series, as volumes 644, 684, 807, 937, 1075, 1264, 1448, 1645, 1848, 2089, 2373, 2676, 3109, 3537, 4009, 4580, 5029, 5577, 6129, 6661, 7354, and 7922.

We thank Russian Foundation for Basic Research, National Research University Higher School of Economics, and Yandex for their financial support.

March 2014

Alexander S. Kulikov
Sergei Kuznetsov
Pavel Pevzner

Organization

Program Committee

Max Alekseyev	George Washington University, USA
Amihood Amir	Bar-Ilan University, Israel
Maxim Babenko	Moscow State University, Russia
Martin Farach-Colton	Rutgers University, USA
Paolo Ferragina	University of Pisa, Italy
Johannes Fischer	Technical University of Dortmund, Germany
Dan Gusfield	University of California at Davis, USA
Roman Kolpakov	Moscow State University, Russia
Gregory Kucherov	Université Paris-Est Marne-la-Vallée, France
Alexander S. Kulikov	St. Petersburg Department of Steklov Institute of Mathematics, Russia (Co-chair)
Juha Kärkkäinen	University of Helsinki, Finland
Gad M. Landau	University of Haifa, Israel
Stefano Lonardi	University of California at Riverside, USA
Ian Munro	University of Waterloo, Canada
S. Muthukrishnan	Rutgers University, USA
Gonzalo Navarro	University of Chile, Chile
Kunsoo Park	Seoul National University, South Korea
Pavel Pevzner	University of California San Diego, USA (Co-chair)
Nadia Pisanti	University of Pisa, Italy
Mikhail Roytberg	Higher School of Economics, Russia
Tatiana Starikovskaya	Higher School of Economics, Russia
Jim Storer	Brandeis University, USA
Jens Stoye	University of Bielefeld, Germany
Esko Ukkonen	University of Helsinki, Finland

Steering Committee

Alberto Apostolico	Georgia Institute of Technology, USA, and IASI-CNR, Italy
Maxime Crochemore	Université Paris-Est, France, and Kings College London, UK
Zvi Galil	Georgia Institute of Technology, USA

Organizing Committee

Stepan Artamonov	Moscow State University, Russia
Maxim Babenko	Moscow State University, Russia
Dmitry Ignatov	Higher School of Economics, Russia
Dmitry Ilvovsky	Higher School of Economics, Russia
Alexander S. Kulikov	St. Petersburg Department of Steklov Institute of Mathematics, Russia
Sergei Kuznetsov	Higher School of Economics, Russia, (Chair)
Dmitry Morozov	Higher School of Economics, Russia
Kamil Salihov	Moscow State University, Russia
Ruslan Savchenko	Moscow State University, Russia
Tatiana Starikovskaya	Higher School of Economics, Russia

Additional Reviewers

Aganezov, Sergey	Kolesnichenko, Ignat
Amit, Mika	Kopczynski, Dominik
Antipov, Dmitry	Labarre, Anthony
Artamonov, Stepan	Levy, Avivit
Bankevich, Anton	Matsieva, Julia
Belazzougui, Djamel	Melsted, Pall
Bowe, Alexander	Mirebrahim, Seyed Hamid
Braga, Marilia	Moret, Bernard
Butman, Ayelet	Nekrich, Yakov
Chikhi, Rayan	Nielsen, Jesper Sindahl
Cicalese, Ferdinando	Ottaviano, Giuseppe
Cording, Patrick Hagge	Paparo, Omer
Cunial, Fabio	Pizzi, Cinzia
Fedorov, Sergey	Polishko, Anton
Feijao, Pedro	Prencipe, Giuseppe
Fertin, Guillaume	Puglisi, Simon
Frid, Yelena	Pérez-Lantero, Pablo
Gagie, Travis	Radoszewski, Jakub
Gawrychowski, Pawel	Rogulenko, Sergey
Giaquinta, Emanuele	Russo, Luis M. S.
Gysel, Rob	Salikhov, Kamil
He, Meng	Savchenko, Ruslan
Hu, Fei	Sirotkin, Alexander
I, Tomohiro	Sirén, Jouni
Inenaga, Shunsuke	St. John, Katherine
Jahn, Katharina	Stevens, Kristian
Jiang, Shuai	Tarasov, Pavel
Kempa, Dominik	Thankachan, Sharma

Tomescu, Alexandru I.
Vasilevskaya, Maria
Venturini, Rossano
Vildhøj, Hjalte Wedel

Vind, Søren
Välimäki, Niko
Willing, Eyla
Wittler, Roland

Abstracts of Invited Talks

Sequence Comparison in the Time of the Deluge

Alberto Apostolico

Georgia Institute of Technology, USA and IASI-CNR, Italy

It is almost fifty years since the appearance of the famous paper entitled “Binary codes capable of correcting deletions, insertions, and reversals” (1966 English translation, Soviet Physics Doklady) by which Vladimir Levenshtein introduced his eponymous measure of string distance. The notion was since to be re-discovered, variously dithered and more or less efficiently computed in distant domains of application. In particular, it provided the platform for half a century of analysis, alignment, search, taxonomy and phylogeny of molecular sequences. As the size and multiplicity of the sequences produced expand to an unprecedented scale, many elegant techniques of the past no longer work. In molecular taxonomy and phylogeny, for instance, the alignment of whole genomes proves both computationally unbearable and hardly significant. In metagenomics, the elucidation of microbiome compositions under conditions of noisy and largely unidentified reference sequences faces steep barriers during assembly and assignment. In recent studies, classical notions are increasingly being complemented or even supplanted by global similarity measures that refer, implicitly or explicitly, to the composition of sequences in terms of their constituent patterns. Such measures hinge more or less uniformly on an underlying notion of relative compressibility, whereby two sequences are similar if either one can be described using mostly pieces from the other. They can free from chores of alignment and assembly. Their computation poses interesting and variously affordable algorithmic problems. This talk will review some such measures, their computation, and their applications.

Repeats in Strings

Maxime Crochemore

King's College London, UK, and Université Paris-Est, France

Large amounts of text are generated every day in the cyberspace via Web sites, emails, social networks, and other communication networks. These text streams need to be analysed to detect critical events or the monitor business for example. An important characteristics to take into account in this setting is the existence of repetitions in texts. Their study constitutes a fundamental area of combinatorics on words due to major applications to string algorithms, data compression, music analysis, and biological sequences analysis, etc. The talk surveys algorithmic methods used to locate repetitive segments in strings. It discusses the notion of runs that encompasses various types of periodicities considered by different authors, as well as the notion of maximal-exponent factors that captures the most significant repeats occurring in a string. The design and analysis of repeat finders rely on combinatorial properties of words and raise a series of open problems in combinatorics on words.

“Stringology” is 30 Years Old

Zvi Galil

Georgia Institute of Technology, USA

This year we have two anniversaries: This is the 25th CPM. CPM1, the first one, was held in Paris in 1990. But in June 1984 there was a sort of a predecessor of CPM1 which can be considered as CPM0. It was a Nato Workshop on combinatorial algorithms on words in Maratea, Italy. The first paper was my paper: “Open Problems in Stringology”. Thus, the term Stringology is exactly 30 years old. The field of stringology is much older. Even though the word stringology cannot be found in any respectable dictionary, Google search yields 23,800 hits; in addition 1640 articles use it, 123 in their title.

In celebration of this 30th anniversary, we will review the status of the 13 open problems of the paper and will introduce some new ones.

How to Think Big

Udi Manber

Google Inc.

“If you have to ask, you’ll never know”
Louis Armstrong (answering a question
about the meaning of “swing”)

In business, making more impact used to mean making more money. But that changed. In 2010 Twitter’s revenues were 40x less than cat litter, but its impact was enormous. Same with Amazon, Google, and Facebook at their beginnings. The key to all of them was that they were doing completely new things. No one thought there was a need to share 140 characters before Twitter, but it turned out to be extremely useful. Very few people thought at the time that shopping on the web, highly relevant search, or communication between friends were important business needs.

When the telephone was introduced to telegraph companies they discounted it. They had a good reason. More than a century later AT&T discounted voice over IP for similar reasons. Yet the 100+ years old AT&T was sold for less than the 5 years old WhatsApp. Things change rather quickly nowadays.

There are some common insights from these stories that can be applied to academic research and I will try to highlight them in this talk.

What's Behind Blast

Gene Myers

Max Planck Institute for Molecular Cell Biology and Genetics

While Blast is one of the most widely used search engines for molecular biology, and while there is a general understanding of how it works, few know the story of how it came about and the theoretical algorithmic result from which it was derived. I will tell the story and explain the theoretical algorithm — an $O(DN^{\text{pow}(D/P)} \log N)$ expected-time algorithm for finding all matches to a query of length P with not more than D differences in a database of length $N \gg P$. Surprisingly, this result, published in early 1994, has to my knowledge not been improved upon to this day.

Table of Contents

On the Efficiency of the Hamming C-Centerstring Problems	1
<i>Amihood Amir, Jessica Fidler, Liam Roditty, and Oren Sar Shalom</i>	
Dictionary Matching with One Gap	11
<i>Amihood Amir, Avivit Levy, Ely Porat, and B. Riva Shalom</i>	
Approximate On-line Palindrome Recognition, and Applications	21
<i>Amihood Amir and Benny Porat</i>	
Computing Minimal and Maximal Suffixes of a Substring Revisited.....	30
<i>Maxim Babenko, Paweł Gawrychowski, Tomasz Kociumaka, and Tatiana Starikovskaya</i>	
Compressed Subsequence Matching and Packed Tree Coloring	40
<i>Philip Bille, Patrick Hagge Cording, and Inge Li Gørtz</i>	
Reversal Distances for Strings with Few Blocks or Small Alphabets.....	50
<i>Laurent Bulteau, Guillaume Fertin, and Christian Komusiewicz</i>	
On Combinatorial Generation of Prefix Normal Words	60
<i>Péter Burcsi, Gabriele Fici, Zsuzsanna Lipták, Frank Ruskey, and Joe Sawada</i>	
Permuted Scaled Matching	70
<i>Ayelet Butman, Noa Lewenstein, and J. Ian Munro</i>	
The Worst Case Complexity of Maximum Parsimony	79
<i>Amir Carmel, Noa Musa-Lempel, Dekel Tsur, and Michal Ziv-Ukelson</i>	
From Indexing Data Structures to de Bruijn Graphs.....	89
<i>Bastien Cazaux, Thierry Lecroq, and Eric Rivals</i>	
Randomized and Parameterized Algorithms for the Closest String Problem	100
<i>Zhi-Zhong Chen, Bin Ma, and Lusheng Wang</i>	
Indexed Geometric Jumbled Pattern Matching.....	110
<i>Stephane Durocher, Robert Fraser, Travis Gagie, Debajyoti Mondal, Matthew Skala, and Sharma V. Thankachan</i>	

An Improved Query Time for Succinct Dynamic Dictionary Matching	120
<i>Guy Feigenblat, Ely Porat, and Ariel Shiftan</i>	
Order-Preserving Pattern Matching with k Mismatches	130
<i>Paweł Gawrychowski and Przemysław Uznański</i>	
Parameterized Complexity Analysis for the Closest String with Wildcards Problem.....	140
<i>Danny Hermelin and Liat Rozenberg</i>	
Computing Palindromic Factorizations and Palindromic Covers On-line	150
<i>Tomohiro I, Shiho Sugimoto, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda</i>	
Compactness-Preserving Mapping on Trees	162
<i>Jan Baumbach, Jiong Guo, and Rashid Ibragimov</i>	
Shortest Unique Substring Query Revisited	172
<i>Atalay Mert İleri, M. Oğuzhan Külekci, and Bojian Xu</i>	
A <i>really</i> Simple Approximation of Smallest Grammar.....	182
<i>Artur Jeż</i>	
Efficient Algorithms for Shortest Partial Seeds in Words	192
<i>Tomasz Kociumaka, Solon P. Pissis, Jakub Radoszewski, Wojciech Rytter, and Tomasz Walen</i>	
Computing k -th Lyndon Word and Decoding Lexicographically Minimal de Bruijn Sequence	202
<i>Tomasz Kociumaka, Jakub Radoszewski, and Wojciech Rytter</i>	
Searching of Gapped Repeats and Subrepetitions in a Word	212
<i>Roman Kolpakov, Mikhail Podolskiy, Mikhail Posypkin, and Nikolay Khrapov</i>	
Approximate String Matching Using a Bidirectional Index	222
<i>Gregory Kucherov, Kamil Salikhov, and Dekel Tsur</i>	
String Range Matching	232
<i>Juha Kärkkäinen, Dominik Kempa, and Simon J. Puglisi</i>	
On Hardness of Several String Indexing Problems	242
<i>Kasper Green Larsen, J. Ian Munro, Jesper Sindahl Nielsen, and Sharma V. Thankachan</i>	
Most Recent Match Queries in On-Line Suffix Trees	252
<i>N. Jesper Larsson</i>	

Encodings for Range Majority Queries	262
<i>Gonzalo Navarro and Sharma V. Thankachan</i>	
On the DCJ Median Problem	273
<i>Mingfu Shao and Bernard M.E. Moret</i>	
Author Index	283