# Lecture Notes in Computer Science 8199

Nieves Brisaboa   Oscar Pedreira
Pavel Zezula (Eds.)

# Similarity Search and Applications

6th International Conference, SISAP 2013
A Coruña, Spain, October 2-4, 2013
Proceedings

Springer

Volume Editors

Nieves Brisaboa
Oscar Pedreira
Universidade da Coruña
Department of Computer Science, A Coruña, Spain
E-mail: {brisaboa, opedreira}@udc.es

Pavel Zezula
Masaryk University
Department of Computer Systems and Communications
Brno, Czech Republic
E-mail: zezula@fi.muni.cz

# Preface

This volume contains the papers presented at the 6th International Conference on Similarity Search and Applications (SISAP 2013), held at A Coruña, Spain, during October 2–4, 2013.

The International Conference on Similarity Search and Applications (SISAP) is an annual forum for researchers and application developers in the area of similarity data management. It aims at the technological problems shared by many application domains, such as data mining, information retrieval, computer vision, pattern recognition, computational biology, geography, biometrics, machine learning, and many others that need similarity searching as a necessary supporting service.

Traditionally, SISAP conferences have put emphasis on the distance-based searching, but in general the conference concerns both the effectiveness and efficiency aspects of any similarity search approach.

In this edition, we had the ambition of widening SISAP's scope to cover even more aspects related to similarity. As a result, we have attracted a higher number of applications covering a wider range of approaches and application domains than in previous editions. In this way, we have achieved a more competitive conference, attractive to a larger part of the computer science community. In the future, we will keep that intent wishing SISAP to become a relevant conference on similarity search, bringing together a wide community of researchers and practitioners and welcoming contributions that range from theoretical aspects to innovative developments for which similarity search plays the central role.

The call for papers welcomed three types of contributions: (*i*) research papers (full or short papers) presenting previously unpublished research contributions, (*ii*) case studies and application papers (short papers) describing existing applications of similarity search in real scenarios, and (*iii*) demo papers describing real or prototype systems for which similarity search technology is a core component, presented at the conference with a system demonstration.

We received 44 submissions, from which 31 were full papers, 11 were short papers, and 2 were demo papers. The Program Committee (PC) comprised 29 researchers from 17 different countries. Each submission was assigned to at least three PC members. Reviews were discussed by the chairs and PC members when the reviews diverged and no sound decision had been reached. The final selection of papers was made by the PC chairs based on the reviews received by each submission. Finally, the conference program includes 19 full papers, 6 short papers, and 2 demo papers, which results in a 61,29% acceptance ratio for full papers and a 54,54% acceptance ration for short papers.

The conference program and the proceedings are organized into five parts. The first one comprises papers dealing with new scenarios or presenting new approaches to similarity search. A second part is devoted to papers proposing

improvements to different methods and techniques for similarity search. The third part focuses on particular metrics and their effectiveness. The fourth part of the conference program includes papers devoted to solutions for similarity search in specific application domains, such as recommender systems, search engines, computational biology, and image and video retrieval. Finally, the last part comprises those papers dealing with efficient implementation and engineering solutions for similarity search in real settings.

The conference program also includes two invited talks from well-known researchers in the field. The first one, "Similarity in Web Search", by Ricardo Baeza-Yates, surveys different aspects of Web search in which similarity search plays an important role, considering the variety of objects that need to be compared, and the nature and features of the metrics involved in each case. The second one, "Large Scale Visual Object Retrieval", by Jiri Matas, presents the state of the art in visual retrieval of specific objects, and describes two new methods for large scale object retrieval.

As in previous editions, the proceedings are published by Springer-Verlag, in the *Lecture Notes in Computer Science* series. A selection of the best papers presented at the conference were recommended for publication in the journal *Information Systems*. The selection of best papers was made by the PC, based on the reviews received by each paper, and on the discussion during the conference.

SISAP conferences are organized by the SISAP initiative (`www.sisap.org`), which aims to become a forum to exchange real-world, challenging and innovative examples of applications, new indexing techniques, common test-beds and benchmarks, source code and up-to-date literature through its web page, serving the similarity search community.

We would like to acknowledge the generous collaboration and financial support from University of A Coruña, Spain (hosting instution), the Fields Institute for Research in Mathematical Sciences, Canada, and the Center for Research and Development in Information Technologies (CITIC) of University of A Coruña. We want to express our gratitude to the PC members for their effort and contribution to the conference. All the submission, reviewing, and proceedings generation processes were carried out through the EasyChair platform.

October 2013                                                    Nieves Brisaboa
                                                                Oscar Pedreira
                                                                  Pavel Zezula

# Organization

## Program Committee Chairs

Nieves R. Brisaboa      Universidade da Coruña, Spain
Pavel Zezula      Masaryk University, Czech Republic

## Program Committee Members

| | |
|---|---|
| Giuseppe Amato | Istituto di Scienza e Tecnologie dellInformazione (CNR), Italy |
| Laurent Amsaleg | Institut de Recherche en Informatique et Systèmes Aléatoires, France |
| Nieves Brisaboa | Universidade da Coruña, Spain |
| Benjamin Bustos | Universidad de Chile, Chile |
| Edgar Chavez | Universidad Nacional Autónoma de México, Mexico |
| Paolo Ciaccia | University of Bologna, Italy |
| Richard Connor | University of Strathclyde, UK |
| Andrea Esuli | Instituto di Scienza e Tecnologie dell'Informazione (CNR), Italy |
| Rosalba Giugno | University of Catania, Italy |
| Michael Houle | National Institute of Informatics, Japan |
| Alexis Joly | Inria, France |
| Björn Jónsson | Reykjavík University, Iceland |
| Daniel Keim | Universität Konstanz, Germany |
| Eamonn Keogh | University of California at Riverside, USA |
| Magnus Lie Hetland | Norwegian University of Science and Technology (NTNU), Norway |
| Yannis Manolopoulos | Aristotle University of Thessaloniki, Greece |
| Rui Mao | Shenzhen University, China |
| Luisa Micó | Universidad de Alicante, Spain |
| Henning Müller | University of Applied Sciences Western Switzerland, Switzerland |
| Gonzalo Navarro | Universidad de Chile, Chile |
| Arlindo Oliveira | Lisbon Technical University, Portugal |
| José Oncina | Universidad de Alicante, Spain |
| Apostolos Papadopoulos | Aristotle University of Thessaloniki, Greece |
| Marco Patella | University of Bologna, Italy |
| Oscar Pedreira | Universidade da Coruña, Spain |
| Vladimir Pestov | University of Ottawa, Canada |

| | |
|---|---|
| Matthias Renz | Ludwig-Maximilians-Universität München, Germany |
| Hanan Samet | University of Maryland, USA |
| Tomas Skopal | Charles University in Prague, Czech Republic |
| Pavel Zezula | Masaryk University, Czech Republic |

## Additional Reviewers

| | |
|---|---|
| Bartolini, Ilaria | Pedreira, Oscar |
| Buisson, Olivier | Reyes, Nora |
| Falchi, Fabrizio | Spretke, David |
| Hoksza, David | Stoffel, Andreas |
| Krulis, Martin | Stoffel, Florian |
| Li, Hao | Symeonidis, Panagiotis |
| Lokoc, Jakub | Tellez, Eric Sadit |
| Moss, Robert | Wanner, Franz |
| Paredes, Rodrigo | |

# Table of Contents

## Metrics and Evaluation

## Applications and Specific Domains

## Implementation and Engineering Solutions

## Demo Papers