# Lecture Notes in Artificial Intelligence 4201

Yasubumi Sakakibara   Satoshi Kobayashi
Kengo Sato   Tetsuro Nishino
Etsuji Tomita (Eds.)

# Grammatical Inference: Algorithms and Applications

8th International Colloquium, ICGI 2006
Tokyo, Japan, September 20-22, 2006
Proceedings

Springer

# Preface

The 8th International Colloquium on Grammatical Inference (ICGI 2006) was held at the University of Electro-Communications (UEC), Tokyo, Japan on September 20-22, 2006. ICGI 2006 was the eighth in a series of successful biennial international conferences in the area of grammatical inference. Previous meetings were held in Essex, UK; Alicante, Spain; Montpellier, France; Ames, Iowa, USA; Lisbon, Portugal; Amsterdam, Netherlands; Athens, Greece. ICGI 2006 was the first conference in this series to be held in Asia. This series of conferences seeks to provide a forum for presentation and discussion of original research papers on all aspects of grammatical inference.

Grammatical inference, the study of learning grammars from data, is an established research field in artificial intelligence, dating back to the 1960s and has been extensively addressed by researchers in automata theory, language acquisition, computational linguistics, machine learning, pattern recognition, computational learning theory and neural networks. ICGI 2006 successively emphasized on the multi-disciplinary nature of the research field and the diverse domains in which grammatical inference is being applied, such as natural language acquisition, computational biology, structural pattern recognition, information retrieval, Web mining, text processing, data compression and adaptive intelligent agents.

We received 44 high-quality papers from 14 countries around the world. The papers were reviewed by three reviewers. Based on the positive comments of the reviewers, 25 full papers were accepted. In addition, we decided to accept 8 short papers for poster presentation. Short papers appear as two-page extended abstracts in a separate section of this volume. The topics of the accepted papers vary from theoretical results of learning algorithms to innovative applications of grammatical inference, and from learning several interesting classes of formal grammars to applications to natural language processing.

In parallel to the submission and reviewing of research papers, a machine translation competition, named Tenjinno, took place. In a separate paper in this volume, the organizers of the competition report on the peculiarities of such an endeavor and some interesting theoretical findings to which they have been led. Last but not least, we were honored by the contributions of our two invited speakers, Yuji Matsumoto, from Nara Institute of Science and Technology, Japan, and Jean-Philippe Vert, from Ecole des Mines de Paris, France. Both invited speakers provided interesting talks on the topics of natural language processing and bioinformatics, and we hope both talks invoked potential applications of grammatical inference.

The editors would like to acknowledge the contribution of the conference's Program Committee and the Additional Reviewers in reviewing the submitted papers and thank the Organizing Committee for their invaluable help in

organizing the conference. Particularly, we would like to thank Colin de la Higuera, Menno van Zaannen, Bradford Starkie, and Dominique Estival for their additional voluntary service to the grammatical inference community, through this conference. We would also like to acknowledge the use of the Cyberchair software, from Borbala online conference services, in the submission and reviewing process. Finally, we are grateful for the generous support and sponsorship of the conference by the University of Electro-Communications, the PASCAL, Inoue foundation for Science, SIG Mathematical Modeling and Problem Solving in Information Processing Society of Japan and New Horizons in Computing (NHC) (Scientific Research on Priority Areas, supported by MEXT, Japan).

September 2006                                          Yasubumi Sakakibara
                                                        Satoshi Kobayashi
                                                           Kengo Sato
                                                        Tetsuro Nishino
                                                          Etsuji Tomita

# Organization

## Conference Chair

Etsuji Tomita                  University of Electro-Communications, Japan

## Technical Program Committee Co-chairs

Yasubumi Sakakibara       Keio University, Japan
Satoshi Kobayashi         University of Electro-Communications, Japan

## Technical Program Committee

| | |
|---|---|
| Naoki Abe | IBM Thomas J. Watson Research Center, USA |
| Pieter Adriaans | Perot Systems Corporation/University of Amsterdam, Netherlands |
| Dana Angluin | Yale University, USA |
| Hiroki Arimura | Hokkaido University, Japan |
| Mitra Basu | City University of New York, USA |
| François Coste | Symbiose, INRIA/IRISA, France |
| Pierre Dupont | University of Louvain, Belgium |
| Henning Fernau | University of Hertfordshire, UK |
| Colin de la Higuera | EURISE, Univ. de St. Etienne, France |
| Vasant Honavar | Iowa State University, USA |
| Chih-Jen Lin | National Taiwan University, Taiwan |
| Laurent Miclet | ENSSAT, Lannion, France |
| Gopalakrishnaswamy Nagaraja | Indian Institute of Technology, India |
| Katsuhiko Nakamura | Tokyo Denki University, Japan |
| Jacques Nicolas | IRISA, France |
| Tim Oates | University of Maryland Baltimore County, USA |
| Arlindo Oliveira | Lisbon Technical University, Portugal |
| Jose Oncina Carratala | Universidad de Alicante, Spain |
| Georgios Paliouras | Inst. of Informatics and Telecommunications, NCSR , Greece |
| Rajesh Parekh | Yahoo!, USA |
| Kengo Sato | CBRC, NAIST, Japan |
| Giora Slutzki | Iowa State University, USA |
| Bradford Starkie | Starkie Enterprise, Australia |
| Eiji Takimoto | Tohoku University, Japan |
| Menno van Zaanen | Universiteit van Amsterdam, Netherlands |
| Enrique Vidal | Universidad Politecnica de Valencia, Spain |
| Osamu Watanabe | Tokyo Institute of Technology, Japan |
| Thomas Zeugmann | Hokkaido University, Japan |

## Additional Reviewers

| | | |
|---|---|---|
| T. Armstrong | J.-C. Janodet | J. Poland |
| L. Becerra-Bonache | H.-U. Krieger | J. M. Vilar |
| M. Bugalho | J. A. Laxminarayana | |
| D. Eisenstat | A. Martins | |

## Organizing Committee Chair

Tetsuro Nishino        University of Electro-Communications, Japan

## Organizing Committee

| | |
|---|---|
| Colin de la Higuera | EURISE, Univ. de St. Etienne, France |
| Kazuhiro Hotta | University of Electro-Communications, Japan |
| Satoshi Kobayashi | University of Electro-Communications, Japan |
| Yoichi Motomura | National Institute of Advanced Industrial Science and Technology, Japan |
| Katsuhiko Nakamura | Tokyo Denki University, Japan |
| Seiya Okubo | University of Electro-Communications, Japan |
| Yasuhiro Tajima | Tokyo University of Agriculture and Technology, Japan |
| Haruhisa Takahashi | University of Electro-Communications, Japan |
| Jun Tarui | University of Electro-Communications, Japan |
| Mitsuo Wakatsuki | University of Electro-Communications, Japan |

## Sponsoring Institutions

University of Electro-Communications

PASCAL Network

Inoue Foundation for Science

SIG Mathematical Modeling and Problem Solving in Information Processing Society of Japan

New Horizons in Computing (NHC) (Scientific Research on Priority Areas, supported by MEXT, Japan)

# Table of Contents

## Invited Papers

## Regular Papers

## Poster Papers