Springer
*Berlin*
*Heidelberg*
*New York*
*Barcelona*
*Hong Kong*
*London*
*Milan*
*Paris*
*Singapore*
*Tokyo*

Arlindo L. Oliveira (Ed.)

# Grammatical Inference: Algorithms and Applications

5th International Colloquium, ICGI 2000
Lisbon, Portugal, September 11-13, 2000
Proceedings

Springer

# Preface

The Fifth International Colloquium on Grammatical Inference (ICGI-2000) was held in Lisbon on September 11–13th, 2000. ICGI-2000 was the fifth in a series of successful biennial international conferences in the area of grammatical inference. Previous conferences were held in Essex, U.K.; Alicante, Spain; Montpellier, France; and Ames, Iowa, USA.

This series of meetings seeks to provide a forum for the presentation and discussion of original research on all aspects of grammatical inference. Grammatical inference, the process of inferring grammar from given data, is a field that is not only challenging from a purely scientific standpoint but also finds many applications in real world problems.

Despite the fact that grammatical inference addresses problems in a relatively narrow area, it uses techniques from many domains, and intersects a number of different disciplines. Researchers in grammatical inference come from fields as diverse as machine learning, theoretical computer science, computational linguistics, pattern recognition and artificial neural networks.

From a practical standpoint, applications in areas such as natural language acquisition, computational biology, structural pattern recognition, information retrieval, text processing and adaptive intelligent agents have either been demonstrated or proposed in the literature.

ICGI-2000 was held jointly with CoNLL-2000, the Computational Natural Language Learning Workshop and LLL-2000, the Second Learning Language in Logic Workshop. The technical program included the presentation of 24 accepted papers (out of 35 submitted) as well as joint sessions with CoNLL and LLL. A tutorial program organized by Gabriel Pereira Lopes took place after the meetings and included tutorials by Raymond Mooney, Gregory Grefenstette, Walter Daelemans, António Ribeiro, Joaquim Ferreira da Silva, Gael Dias, Nuno Marques, Vitor Rossio, João Balsa and Alexandre Agostini. The joint realization of these events represents a unique opportunity for researchers in these related fields to interact and exchange ideas.

I would like to thank Claire Nédellec, Claire Cardie, Walter Daelemans, Colin de la Higuera and Vasant Honavar for their help in several aspects of the organization; the members of the technical program committee and the reviewers for their careful evaluation of the submissions; the members of the local organizing committee, Ana Teresa Freitas and Ana Fred, for their help in setting up the event; and Ana de Jesus for her invaluable secretarial support.

September 2000                                        Arlindo Oliveira
                                                Technical Program Chair

## Technical Program Committee

| | |
|---|---|
| Pieter Adriaans | Syllogic/University of Amsterdam, The Netherlands |
| Michael Brent | Johns Hopkins University, USA |
| Walter Daelemans | Tilburg University, The Netherlands |
| Pierre Dupont | University de St. Etienne, France |
| Dominique Estival | Syrinx Speech Systems, Australia |
| Ana Fred | Lisbon Technical University, Portugal |
| Jerry Feldman | ICSI, Berkeley, USA |
| Lee Giles | NEC Research Institute, USA |
| Colin de la Higuera | EURISE, University de St. Etienne, France |
| Vasant Honavar | Iowa State University, USA |
| Laurent Miclet | ENSSAT, France |
| G. Nagaraja | Indian Institute of Technology, India |
| Jacques Nicolas | IRISA, France |
| Arlindo Oliveira | INESC/IST, Portugal |
| Jose Oncina Carratala | Universidade de Alicante, Spain |
| Rajesh Parekh | Allstate Research and Planning Center, USA |
| Lenny Pitt | University of Illinois at Urbana-Champaign, USA |
| Yasubumi Sakakibara | Tokyo Denki University, Japan |
| Arun Sharma | University of New South Wales, Australia |
| Giora Slutzki | Iowa State University, USA |
| Esko Ukkonen | University of Helsinki, Finland |
| Stefan Wermter | University of Sunderland, UK |
| Enrique Vidal | University Politecnica de Valencia, Spain |
| Thomas Zeugmann | Kyushu University, Japan |

## Organizing Committee

| | |
|---|---|
| Conference Chair: | Arlindo Oliveira, INESC/IST |
| Tutorials: | Gabriel Pereira Lopes, Universidade Nova de Lisboa |
| Local Arrangements: | Ana Fred, Lisbon Technical University |
| Social Program: | Ana Teresa Freitas, INESC/IST |
| Secretariat: | Ana de Jesus, INESC |

## Additional Reviewers

Daniel Gildea
Mitch Harris
Satoshi Kobayashi
Eric Martin
Franck Thollard
Takashi Yokomori

# Table of Contents