

# Lecture Notes in Artificial Intelligence

1433

Subseries of Lecture Notes in Computer Science

Edited by J. G. Carbonell and J. Siekmann

Lecture Notes in Computer Science

Edited by G. Goos, J. Hartmanis and J. van Leeuwen

Vasant Honavar Giora Slutzki (Eds.)

# Grammatical Inference

4th International Colloquium, ICGI-98  
Ames, Iowa, USA, July 12-14, 1998  
Proceedings



Springer

Series Editors

Jaime G. Carbonell, Carnegie Mellon University, Pittsburgh, PA, USA  
Jörg Siekmann, University of Saarland, Saarbrücken, Germany

Volume Editors

Vasant Honavar

Giora Slutzki

Iowa State University, Department of Computer Science

226 Atanasoff Hall, Ames, Iowa, 50011-1040, USA

E-mail: {honavar,slutzki}@cs.iastate.edu

Cataloging-in-Publication Data applied for

**Die Deutsche Bibliothek - CIP-Einheitsaufnahme**

**Grammatical inference : 4th international colloquium ; proceedings / ICGI-98, Ames, Iowa, USA, July 12 - 14, 1998. Vasant Honavar ; Giora Slutzki (ed.). - Berlin ; Heidelberg ; New York ; Barcelona ; Budapest ; Hong Kong ; London ; Milan ; Paris ; Singapore ; Tokyo : Springer, 1998**

**(Lecture notes in computer science ; Vol. 1433 : Lecture notes in artificial intelligence)**

**ISBN 3-540-64776-7**

CR Subject Classification (1991): I.2, F.4.2-3, I.5.1, I.5.4, J.5

ISBN 3-540-64776-7 Springer-Verlag Berlin Heidelberg New York

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, re-use of illustrations, recitation, broadcasting, reproduction on microfilms or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permission for use must always be obtained from Springer-Verlag. Violations are liable for prosecution under the German Copyright Law.

© Springer-Verlag Berlin Heidelberg 1998  
Printed in Germany

Typesetting: Camera ready by author

SPIN 10637663 06/3142 - 5 4 3 2 1 0 Printed on acid-free paper

## Preface

Grammatical Inference, also referred to as automata induction, grammar induction, and automatic language acquisition, is the process of learning of grammars, automata, and languages from data (examples, queries, etc.). Machine learning of grammars finds a variety of applications in syntactic pattern recognition, adaptive intelligent agents, diagnosis, computational biology, systems modelling, prediction, natural language acquisition, data mining and knowledge discovery.

Historically, grammatical inference has been studied by researchers in several research communities including: Information Theory, Formal Languages, Automata Theory, Language Acquisition, Computational Linguistics, Machine Learning, Pattern Recognition, Computational Learning Theory, Neural Networks, etc. These different communities have been largely isolated.

Perhaps one of the first attempts to bring together researchers working on grammatical inference for an interdisciplinary exchange of research results took place under the aegis of the First Colloquium on Grammatical Inference held at the University of Essex in United Kingdom in April 1993. This was followed by the (second) International Colloquium on Grammatical Inference, held at Alicante in Spain and the Third International Colloquium on Grammatical Inference, held at Montpellier in France. Following the success of these events and the Workshop on Automata Induction, Grammatical Inference, and Language Acquisition, held in conjunction with the International Conference on Machine Learning at Nashville in United States in July 1997, it was deemed appropriate to hold the Fourth International Colloquium on Grammatical Inference (ICGI '98) in United States.

ICGI '98 consisted of both refereed papers as well as two invited papers and an invited tutorial. Approximately 35 papers were received for review from Europe, United States, India, Japan, and Australia. Each submitted paper was reviewed for technical soundness, originality, clarity of presentation, and relevance by at least 2 members of the program committee. A total of 22 papers were accepted for presentation at the conference. The contributed and invited papers cover a wide range of topics in the theory as well as the applications of grammatical inference, automata induction, and language learning.

We are grateful to the members of the Technical Program Committee for reviewing the papers, and the Local Arrangements Committee and Ms. Margie Poorman and Ms. Janet Gardner of the Extended and Continuing Education Division of Iowa State University for their help with the organization of the conference. We would like to thank Professor James Vary (the Director of the Institute for Theoretical and Applied Physics (IITAP)) Professor John Mayfield (the Associate Dean of the Graduate College), and Professor Arthur Oldehoeft (the Chair of Computer Science) at Iowa State University for their support of the conference. We would like to thank Professor Laurent Miclet and Professor Colin de la Higuera for their help with many aspects of the conference. We

are grateful to Professor Jerry Feldman and Professor Alvis Brazma for kindly agreeing to give invited talks. Professor Jack Lutz graciously agreed to give a tutorial on Kolmogorov Complexity and its Applications. We would like to thank the attendees for making ICGI '98 a success. We are grateful to the the editorial staff of Springer-Verlag for putting together the conference proceedings.

July 1998

Vasant Honavar and Giora Slutzki  
Program Chairs  
ICGI '98.

# Organization

ICGI '98 was organized at Iowa State University (ISU) in Ames, Iowa from July 12 through July 14, 1998 by Vasant Honavar and Giora Slutzki with the assistance of the Technical Program Committee, the Local Arrangements Committee, and the staff of ISU Extended and Continuing Education Division.

ICGI '98 was held in cooperation with the American Association for Artificial Intelligence (AAAI), the IEEE Systems, Man and Cybernetics Society, and the Association for Computational Linguistics (ACL) Special Interest Group on Natural Language Learning.

ICGI '98 was cosponsored by the International Institute of Theoretical and Applied Physics (IITAP), the Iowa Computational Biology Laboratory, the Complex Adaptive Systems Group, the Artificial Intelligence Research Laboratory and the Department of Computer Science at Iowa State University.

## Program Committee

### Program Chairs

Vasant Honavar and Giora Slutzki, Iowa State University, USA

### Technical Program Committee

R. Berwick, MIT, USA  
A. Brazma, European Bioinformatics Institute, Cambridge, UK.  
M. Brent, Johns Hopkins University, USA  
C. Cardie, Cornell University, USA  
W. Daelemans, Tilburg University, Netherlands  
D. Dowe, Monash University, Australia  
P. Dupont, Univ. St. Etienne, France.  
D. Estival, University of Melbourne, Australia  
J. Feldman, International Computer Science Institute, Berkeley, USA  
L. Giles, NEC Research Institute, Princeton, USA  
J. Gregor, University of Tennessee, USA  
C. de la Higuera, LIRMM, France  
A. Itai, Technion, Israel  
T. Knuutila, University of Turku, Finland  
J. Koza, Stanford University, USA  
M. Li, University of Waterloo, Canada  
E. Makinen, University of Tampere, Finland  
L. Miclet, ENSSAT, Lannion, France.  
G. Nagaraja, Indian Institute of Technology, Bombay, India  
H. Ney, University of Technology, Aachen, Germany  
J. Nicolas, IRISA, France

- R. Parekh, Allstate Research and Planning Center, USA
- L. Pitt, University of Illinois at Urbana-Champaign, USA
- D. Powers, Flinders University, Australia
- L. Reeker, National Science Foundation, USA
- Y. Sakakibara, Tokyo Denki University, Japan
- C. Samuelsson, Lucent Technologies, USA
- A. Sharma, University of New South Wales, Australia.
- E. Vidal, U. Politecnica de Valencia, Spain

**Local Arrangements Committee**

- Dale Grosvenor, Iowa State University, USA.
- K. Balakrishnan, Iowa State University, USA.
- R. Bhatt, Iowa State University, USA
- J. Yang, Iowa State University, USA.

## Table of Contents

Results of the Abbadingo One DFA Learning Competition and a New Evidence-Driven State Merging Algorithm <i>Kevin J. Lang, Barak A. Pearlmutter, and Rodney A. Price</i>	1
Learning k-Variable Pattern Languages Efficiently Stochastically Finite on Average from Positive Data <i>Peter Rossmanith and Thomas Zeugmann</i>	13
Meaning Helps Learning Syntax <i>Isabelle Tellier</i>	25
A Polynomial Time Incremental Algorithm for Learning DFA <i>Rajesh Parekh, Codrin Nichitiu, and Vasant Honavar</i>	37
The Data Driven Approach Applied to the OSTIA Algorithm <i>José Oncina</i>	50
Grammar Model and Grammar Induction in the System NL PAGE <i>Vlado Kešelj</i>	57
Approximate Learning of Random Subsequential Transducers <i>Antonio Castellanos</i>	67
Learning Stochastic Finite Automata from Experts <i>Colin de la Higuera</i>	79
Learning Deterministic Finite Automaton with a Recurrent Neural Network <i>Laura Firoiu, Tim Oates, and Paul R. Cohen</i>	90
Applying Grammatical Inference in Learning a Language Model for Oral Dialogue <i>Jacques Chodorowski and Laurent Miclet</i>	102
Real Language Learning <i>Jerome A. Feldman</i>	114
A Stochastic Search Approach to Grammar Induction <i>Hugues Juillé and Jordan B. Pollack</i>	126
Transducer-Learning Experiments on Language Understanding <i>David Picó and Enrique Vidal</i>	138
Locally Threshold Testable Languages in Strict Sense: Application to the Inference Problem <i>José Ruiz, Salvador España, and Pedro García</i>	150

Learning a Subclass of Linear Languages from Positive Structural Information <i>José M. Sempere and G. Nagaraja</i>	162
Grammatical Inference in Document Recognition <i>Alexander S. Saidi and Souad Tayeb-bey</i>	175
Stochastic Inference of Regular Tree Languages <i>Rafael C. Carrasco, José Oncina, and Jorge Calera</i>	187
How Considering Incompatible State Mergings May Reduce the DFA Induction Search Tree <i>François Coste and Jacques Nicolas</i>	199
Learning Regular Grammars to Model Musical Style: Comparing Different Coding Schemes <i>Pedro P. Cruz-Alcázar and Enrique Vidal-Ruiz</i>	211
Learning a Subclass of Context-Free Languages <i>J.D. Emerald, K.G. Subramanian, and D.G. Thomas</i>	223
Using Symbol Clustering to Improve Probabilistic Automaton Inference <i>Pierre Dupont and Lin Chase</i>	232
A Performance Evaluation of Automatic Survey Classifiers <i>Peter Viechnicki</i>	244
Pattern Discovery in Biosequences <i>Alvis Brāzma, Inge Jonassen, Jaak Vilo, and Esko Ukkonen</i>	257
<b>Author Index</b>	<b>271</b>