Lecture Notes in Computer Science

Lecture Notes in Artificial Intelligence 14127

Founding Editor Jörg Siekmann

Series Editors

Randy Goebel, University of Alberta, Edmonton, Canada Wolfgang Wahlster, DFKI, Berlin, Germany Zhi-Hua Zhou, Nanjing University, Nanjing, China The series Lecture Notes in Artificial Intelligence (LNAI) was established in 1988 as a topical subseries of LNCS devoted to artificial intelligence.

The series publishes state-of-the-art research results at a high level. As with the LNCS mother series, the mission of the series is to serve the international R & D community by providing an invaluable service, mainly focused on the publication of conference and workshop proceedings and postproceedings.

Davide Calvaresi · Amro Najjar · Andrea Omicini · Reyhan Aydogan · Rachele Carli · Giovanni Ciatto · Yazan Mualla · Kary Främling Editors

Explainable and Transparent AI and Multi-Agent Systems

5th International Workshop, EXTRAAMAS 2023 London, UK, May 29, 2023 Revised Selected Papers



Editors Davide Calvaresi University of Applied Sciences and Arts Western Switzerland Sierre, Switzerland

Andrea Omicini 🖻 Alma Mater Studiorum, Università di Bologna Bologna, Italy

Rachele Carli 💿 Alma Mater Studiorum, Università di Bologna Bologna, Italy

Yazan Mualla Université de Technologie de Belfort-Montbéliard Belfort Cedex, France Amro Najjar Luxembourg Institute of Science and Technology Esch-sur-Alzette, Luxembourg

Reyhan Aydogan D Ozyegin University Istanbul, Türkiye

Giovanni Ciatto D Alma Mater Studiorum, Università di Bologna Bologna, Italy

Kary Främling Umeå University Umeå, Sweden

ISSN 0302-9743 ISSN 1611-3349 (electronic) Lecture Notes in Artificial Intelligence ISBN 978-3-031-40877-9 ISBN 978-3-031-40878-6 (eBook) https://doi.org/10.1007/978-3-031-40878-6

LNCS Sublibrary: SL7 - Artificial Intelligence

© The Editor(s) (if applicable) and The Author(s), under exclusive license to Springer Nature Switzerland AG 2023

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.

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.

The publisher, the authors, and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Switzerland AG The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

Preface

AI research has made several significant breakthroughs that has boosted its adoption in several domains, impacting our lives on a daily basis. Nevertheless, such widespread adoption of AI-based systems has raised concerns about their foreseeability and controllability and led to initiatives to "slow down" AI research. While such debates have mainly taken place in the media, several other research works have emphasized that achieving trustworthy and responsible AI would necessitate making AI more transparent and explainable.

Not only would eXplainable AI (XAI) increase acceptability, avoid failures, and foster trust, but it would also comply with relevant (inter)national regulations and highlight these new technologies' limitations and potential.

In 2023, the fifth edition of the EXplainable and TRAnsparent AI and Multi-Agent Systems (EXTRAAMAS) continued the successful track initiated in 2019 in Montreal and followed by the 2020 to 2022 editions (virtual due to the COVID-19 pandemic circumstances). Finally, EXTRAAMAS 2023 was held in person and proposed bright presentations, a stimulating keynote (titled "Untrustworthy AI" given by Jeremy Pitt, Imperial College London), and engaging discussions and Q&A sessions.

Overall, EXTRAAMAS 2023 welcomed contributions covering areas including (i) XAI in symbolic and subsymbolic AI, (ii) XAI in negotiation and conflict resolution, (iii) Explainable Robots and Practical Applications, and (iv) (X)AI in Law and Ethics.

EXTRAAMAS 2023 received 26 submissions. Each submission underwent a rigorous single-blind peer-review process (three to five reviews per paper). Eventually, 16 papers were accepted and collected in this volume.

Each paper was presented in person (with the authors' consent, they are available on the EXTRAAMAS website¹). Finally, The Main Chairs would like to thank the special track chairs, publicity chairs, and Program Committee for their valuable work, as well as the authors, presenters, and participants for their engagement.

June 2023

Davide Calvaresi Amro Najjar Andrea Omicini Kary Främling

¹ https://extraamas.ehealth.hevs.ch/.

Organization

General Chairs

Davide Calvaresi	University of Applied Sciences and Arts Western Switzerland, Switzerland
Amro Najjar	Luxembourg Institute of Science and Technology,
	Luxembourg
Andrea Omicini	Alma Mater Studiorum, Università di Bologna, Italy
Kary Främling	Umeå University, Sweden

Special Track Chairs

Ozyegin University, Turkiye
Alma Mater Studiorum, Università di Bologna,
Italy
UTBM, France
Alma Mater Studiorum, Università di Bologna, Italy

Publicity Chairs

Yazan Mualla	UTBM, France
Benoit Alcaraz	University of Luxembourg, Luxembourg
Rachele Carli	Alma Mater Studiorum, Università di Bologna,
	Italy

Advisory Board

Tim Miller Michael Schumacher

Virginia Dignum Leon van der Torre University of Melbourne, Australia University of Applied Sciences and Arts Western Switzerland, Switzerland Umeå University, Sweden University of Luxembourg, Luxembourg

Program Committee

Andrea Agiollo Remy Chaput Paolo Serrani Federico Sabbatini Kary Främling Davide Calvaresi Rachele Carli Victor Contreras Francisco Rodríguez Lera Bartłomiej Kucharzyk Timotheus Kampik Igor Tchappi Eskandar Kouicem Mickaël Bettinelli Rvuta Arisaka Alaa Daoud Avleen Malhi Minal Patil Yazan Mualla Takayuki Ito Lora Fanda Marina Paolanti Arianna Rossi Joris Hulstijn Mahjoub Dridi Thiago Raulino Giuseppe Pisano Katsuhide Fujita Jomi Fred Hubner Roberta Calegari Hui Zhao Niccolo Marini Salima Lamsiyah Stephane Galland

Giovanni Ciatto

University of Bologna, Italy University of Lyon 1, France Universià Politecnica delle Marche, Italy University of Bologna, Italy Umeå University, Sweden University of Applied Sciences and Arts Western Switzerland, Switzerland University of Bologna, Italy University of Applied Sciences and Arts Western Switzerland, Switzerland University of León, Spain Jagiellonian University, Poland University of Umeå, Sweden, Signavio GmbH, Germany University of Luxembourg, Luxembourg Université Grenoble Alpes, France Université Savoie Mont Blanc, France Kyoto University, Japan **INSA Rouen-Normandie**, France Bournemouth University, UK Umeå University, Sweden UTBM. France Kyoto University, Japan University of Geneva, Switzerland Università Politecnica delle Marche, Italy University of Luxembourg, Luxembourg University of Luxembourg, Luxembourg UTBM. France Federal University of Santa Catarina, Brasil Univeristy of Bologna, Italy Tokyo University of Agriculture and Technology, Japan Federal University of Santa Catarina, Brazil University of Bologna, Italy Tongji University, China University of Applied Sciences Western Switzerland, Switzerland University of Luxembourg, Luxembourg UTBM, France University of Bologna, Italy

Matteo Magnini Viviana Mascardi Giovanni Sileno Sarath Sreedharan University of Bologna, Italy University of Genoa, Italy University of Amsterdam, The Netherlands Arizona State University, USA

Contents

Explainable Agents and Multi-Agent Systems	
Mining and Validating Belief-Based Agent Explanations Ahmad Alelaimat, Aditya Ghose, and Hoa Khanh Dam	
Evaluating a Mechanism for Explaining BDI Agent Behaviour Michael Winikoff and Galina Sidorenko	18
A General-Purpose Protocol for Multi-agent Based Explanations Giovanni Ciatto, Matteo Magnini, Berk Buzcu, Reyhan Aydoğan, and Andrea Omicini	38
Dialogue Explanations for Rule-Based AI Systems Yifan Xu, Joe Collenette, Louise Dennis, and Clare Dixon	59
Estimating Causal Responsibility for Explaining Autonomous Behavior Saaduddin Mahmud, Samer B. Nashed, Claudia V. Goldman, and Shlomo Zilberstein	78
Explainable Machine Learning	

The Quarrel of Local Post-hoc Explainers for Moral Values Classification in Natural Language Processing Andrea Agiollo, Luciano Cavalcante Siebert, Pradeep Kumar Murukannaiah, and Andrea Omicini	
Bottom-Up and Top-Down Workflows for Hypercube- And Clustering-Based Knowledge Extractors Federico Sabbatini and Roberta Calegari	116
Imperative Action Masking for Safe Exploration in Reinforcement Learning	130
Reinforcement Learning in Cyclic Environmental Changes for Agents in Non-Communicative Environments: A Theoretical Approach <i>Fumito Uwano and Keiki Takadama</i>	143

xii	Contents

Inherently Interpretable Deep Reinforcement Learning Through Online	
Mimicking	160
Andreas Kontogiannis and George A. Vouros	
Counterfactual, Contrastive, and Hierarchical Explanations	
with Contextual Importance and Utility	180
Kary Främling	

Cross-Domain Applied XAI

Explanation Generation via Decompositional Rules Extraction for Head and Neck Cancer Classification Victor Contreras, Andrea Bagante, Niccolò Marini, Michael Schumacher, Vincent Andrearczyk, and Davide Calvaresi	
Metrics for Evaluating Explainable Recommender Systems Joris Hulstijn, Igor Tchappi, Amro Najjar, and Reyhan Aydoğan	212
Leveraging Imperfect Explanations for Plan Recognition Problems Ahmad Alelaimat, Aditya Ghose, and Hoa Khanh Dam	231
Reinterpreting Vulnerability to Tackle Deception in Principles-Based XAI for Human-Computer Interaction	249
Enhancing Wearable Technologies for Dementia Care: A Cognitive Architecture Approach Matija Franklin, David Lagnado, Chulhong Min, Akhil Mathur, and Fahim Kawsar	270
Author Index	281