# The Friendship Field - an Agent-Based Model on Dyadic Friendship Formation Driven by Social Battery

Flaminio Squazzoni   *Editor*

# Advances in Social Simulation

Proceedings of the 17th Social Simulation Conference, European Social Simulation Association

Springer

# Springer Proceedings in Complexity

Springer Proceedings in Complexity publishes proceedings from scholarly meetings on all topics relating to the interdisciplinary studies of complex systems science. Springer welcomes book ideas from authors. The series is indexed in Scopus.

Proposals must include the following:

- name, place and date of the scientific meeting
- a link to the committees (local organization, international advisors etc.)
- scientific description of the meeting
- list of invited/plenary speakers
- an estimate of the planned proceedings book parameters (number of pages/articles, requested number of bulk copies, submission deadline)

Submit your proposals to: Hisako.Niko@springer.com

Flaminio Squazzoni
Editor

# Advances in Social Simulation

Proceedings of the 17th Social Simulation Conference, European Social Simulation Association

Springer

*Editor*
Flaminio Squazzoni
Dipartimento di Scienze Sociali e Politiche
Università degli Studi di Milano
Milan, Italy

# Committees

## SSC2022 Conference Chair

Flaminio Squazzoni, University of Milan, Italy

## SSC2022 Local Organizational Staffs

Federico Bianchi, University of Milan, Italy
Marco Cremonini, University of Milan, Italy
Carlo Debernardi, University of Milan, Italy
Francesco Renzini, University of Milan, Italy

## Programme Committee

Diana Francisca Adamatti, Universidade Federal do Rio Grande, Brazil
Petra Ahrweiler, Johannes Gutenberg University Mainz, Germany
Fred Amblard, IRIT—University Toulouse 1 Capitole, France
Giulia Andrighetto, CNR, Rome, Italy
Luis Antunes, University of Lisbon, Portugal
David Anzola, Universidad del Rosario, Colombia
Peter Barbrook-Johnson, University of Oxford, United Kingdom
Federico Bianchi, University of Milan, Italy
Mike Bithell, University of Cambridge, United Kingdom
Riccardo Boero, Norwegian Institute for Air Research, Norway
Melania Borit, UiT—The Arctic University of Norway, Norway
Dino Carpentras, University of Limerick, Ireland
Emile Chappin, Technical University of Delft, Netherlands

Kevin Chapuis, IRD—Institute of Research for Development, Montpellier, France
Edmund Chattoe-Brown, University of Leicester, United Kingdom
Andrew Crooks, University at Buffalo, United States
Marcin Czupryna, Cracow University of Economics, Poland
Paola D'Orazio, Ruhr University Bochum, Germany
Natalie Davis, Lancaster University, United Kingdom
Tristan de Wildt, Technical University of Delft, Netherlands
Guillaume Deffuant, Cemagref, France
Frank Dignum, Umeå Universitet, Sweden
Alexis Drogoul, IRD—Institute of Research for Development, Vietnam
Bruce Edmonds, Manchester Metropolitan University Business School, United Kingdom
Joshua M. Epstein, New York University College of Global Public Health, United States
Andreas Flache, University of Groningen, Netherlands
Christopher Frantz, Norwegian University of Science and Technology, Norway
Simone Gabbriellini, Manent.AI, Italy
Cesar Garcia-Diaz, Pontificia Universidad Javeriana, Colombia
Amineh Ghorbani, Technical University of Delft, Netherlands
Francesca Giardini, University of Groningen, Netherlands
Nigel Gilbert, University of Surrey, United Kingdom
Francisco Grimaldo, Universitat de València, Spain
Rainer Hegselmann, Bayreuth University, Germany
Sebastian Hoffmann, TU Dortmund University, Germany
Gertjan Hofstede, Wageningen University, Netherlands
Sascha Holzhauer, Universität Kassel, Germany
Sylvie Huet, INRAE, France
Luis R. Izquierdo, University of Burgos, Spain
Wander Jager, University of Groningen, Netherlands
Toshiya Kaihara, Kobe University, Japan
Bogumił Kamiński, SGH Warsaw School of Economics, Poland
Andreas Koch, University of Salzburg, Austria
Friedrich Krebs, University Kassel, Germany
Stephan Leitner, University of Klagenfurt, Austria
Pablo Lucas, University College Dublin, Ireland
Katharina Luckner, University of Hamburg, Germany
Ruth Meyer, Centre for Policy Modelling, United Kingdom
Selcan Mutgan, Linköping University, Sweden
Kavin Preethi Narasimhan, University of Surrey, United Kingdom
Martin Neumann, Johannes Gutenberg University Mainz, Germany
Leila Niamir, IIASA-International Institute for Applied Systems Analysis Vienna, Austria
Isamu Okada, Faculty of Business Administration, Soka University Japan, Japan
Pawel Oleksy, Cracow University of Economics, Poland

Jonathan Ozik, Argonne National Laboratory and The University of Chicago, United States
Bill Rand, North Carolina State University, United States
Jessica Reale, Ruhr Universität Bochum, Germany
Michael Roos, Ruhr Universität Bochum, Germany
Juliette Rouchier, CNRS LAMSADE, France
Jordi Sabater Mir, IIIA-CSIC, Spain
Mauricio Salgado, Centro de Estudios Publicos, Chile
Geeske Scholz, Institute of Environmental Systems Research University of Osnabrück, Germany
Davide Secchi, University of Southern Denmark, Denmark
Roman Seidl, Leibniz University Hannover, Germany
Leron Shults, University of Agder, Norway
Jaime Simão Sichman, Universidade de São Paulo, Brazil
Peer-Olaf Siebers, University of Nottingham, United Kingdom
Małgorzata Snarska, Cracow University of Economics, Poland
Flaminio Squazzoni, University of Milan, Italy
Timo Szczepanska, UiT—The Arctic University of Norway, Norway
Aron Szekely, Collegio Carlo Alberto, Italy
Przemyslaw Szufel, Warsaw School of Economics, Poland
Klaus Troitzsch, University of Koblenz and Landau, Germany
Lois Vanhée, Umeå Universitet, Sweden
Harko Verhagen, Stockholm University, Sweden
Friederike Wall, Alpen-Adria-Universitaet Klagenfurt, Austria
Nanda Wijermans, Stockholm University, Sweden
Hang Xiong, Huazhong Agricultural University, China

## External Reviewers

Morteza Alaeddini, Grenoble Alpes University, France
Gayanga Bandara Herath, University of Southern Denmark, Denmark
Andrew Bell, Boston University, United States
Stefano Benincasa, University of Southern Denmark, Denmark
Dario Blanco-Fernandez, University of Klagenfurt, Austria
Christine Boshuijzen-van Burken, University of New South Wales, Australia
Michele Catalano, International Institute for Applied Systems Analysis, Austria
Carlo Debernardi, University of Milan, Italy
Yue Dou, University of Twente, Netherlands
Siavash Farahbakhsh, Flanders Research Institute for Agriculture, Fisheries and Food, Belgium
Thomas Feliciani, University College Dublin, Ireland
Tatiana Filatova, Delft University of Technology, Netherlands
Andreas Flache, University of Groningen, Netherlands

Christopher Frantz, Norwegian University of Science and Technology, Norway
Dehua Gao, Shandong Technology and Business University, China
Benoit Gaudou, Toulouse 1 Capitole University, France
Jiaqi Ge, University of Leeds, United Kingdom
Yusuke Goto, Shibaura Institute of Technology, Japan
Arvid Horned, Umeå University, Sweden
Aashis Joshi, Delft University of Technology, Netherlands
Christian Kammler, Umeå University, Sweden
Jean-Daniel Kant, Sorbonne University, France
Ravshanbek Khodzhimatov, University of Klagenfurt, Austria
Hajime Kita, Kyoto University, Japan
Setsuya Kurahashi, University of Tsukuba, Japan
Marco Janssen, Arizona State University, United States
Silvia Leoni, INFN, Italy
Katharina Luckner, University of Hamburg, Germany
Gianluca Manzo, Sorbonne University, France
Giovanni Massari, Polytechnic University of Bari, Italy
Patrick Mellacher, University of Graz, Austria
René Mermella, Umeå University, Sweden
Azhar Mohd Ibrahim, International Islamic University Malaysia, Malaysia
Tadahiko Murata, Kansai University, Japan
Reinhard Neck, University of Klagenfurt, Austria
Brayton Noll, Delft University of Technology, Netherlands
Rocco Paolillo, Bremen International Graduate School of Social Sciences, Germany
Paolo Pellizzari, Università Ca' Foscari Venezia, Italy
Jannick Plähn, Hamburg University of Technology, Germany
Eleonora Priori, University of Turin, Italy
Robin Purshouse, University of Sheffield, United Kingdom
Shyaam Ramkumar, University of Milan, Italy
Alexandra Rausch, University of Klagenfurt, Austria
Fernanda Reintgen, University of Groningen, Netherlands
Severin Reissl, European Institute on Economics and the Environment, Italy
Francesco Renzini, University of Milan, Italy
Simone Righi, Università Ca' Foscari Venezia, Italy
Nikitas Sgouros, University of Pisa, Italy
John Stevenson, California Institute of Technology, United States
Emilio Sulis, University of Turin, Italy
Timo Szczepanska, The Arctic University of Norway, Norway
Alessandro Taberna, Delft University of Technology, Netherlands
Shingo Takahashi, Waseda University, Japan
Andrea Teglio, Università Ca' Foscari Venezia, Italy
Pietro Terna, University of Turin, Italy
Marco Tolotti, Università Ca' Foscari Venezia, Italy
Daniel Torren Peraire, Autonomous University of Barcelona, Spain
Lois Vanhee, Umeå University, Sweden

Pierluigi Velucci, Roma Tre University, Italy
Eva Vriens, National Research Council, Italy
Tae-Sub Yun, Korea Advanced Institute of Science and Technology, Korea

# Preface

This book includes the proceedings of the Social Simulation Conference 2022, the 17th annual conference of ESSA—The European Social Simulation Association, held at the University of Milan, Italy, on 12–16 September 2022. Among the various initiatives to promote the development of social simulation research, education and application in Europe, ESSA has organized since 2003 an annual conference that has become the major international annual event for scholars and practitioners interested in the latest developments of this interdisciplinary field of research.

Thanks to the support of the Department of Social and Political Sciences of the University of Milan and the organizational staff of the BehaveLab, the 2022 edition of the conference attracted more than 200 participants in a hybrid setting. The conference included 25 tracks with 145 papers, out of which 46 were selected for publication in these proceedings. The conference papers were reviewed by the programme committee members and a group of external reviewers. The authors of these book chapters have greatly benefited from their feedback and comments and so I would like to take the opportunity here to express my gratitude to them.

Milan, Italy                                                                                   Flaminio Squazzoni

# Contents

## Tools and Methods

# Contributors

**Naphtali Abudarham**  Rafael—Advanced Defense Systems Ltd., Gazit Institute, Tel-Aviv, Israel

**Fabian Adelt**  Technology Studies Group, Faculty of Social Sciences, TU Dortmund University, Dortmund, Germany

**Morteza Alaeddini**  Grenoble Informatics Laboratory (LIG), Université Grenoble Alpes, Grenoble, France

**Aarthi Ananthanarayanan**  Ocean Conservancy, Washington, D.C., USA

**Patrycja Antosz**  NORCE Norwegian Research Centre AS, Bergen, Norway

**Willem Auping**  Faculty of Technology Policy and Management, Delft University of Technology, Jaffalaan 5, 2628 BX Delft, The Netherlands

**Rachel J. Bacon**  Center for Mind and Culture, Boston, MA, USA

**Richard Bailey**  University of Oxford, Oxford, UK

**Matteo Barsanti**  School of Architecture, Civil and Environmental Engineering, EPFL, Ecublens VD, Switzerland

**Lucia Bellora-Bienengräber**  Department of Accounting and Auditing, Faculty of Economics and Business, University of Groningen, AE, Groningen, The Netherlands; Wilbur O. and Ann Powers College of Business, School of Accountancy, Clemson, USA

**Federico Bianchi**  Department of Social and Political Sciences, University of Milan, Milan, Italy

**Claudia Binder**  School of Architecture, Civil and Environmental Engineering, EPFL, Ecublens VD, Switzerland

**Darío Blanco-Fernández**  University of Klagenfurt, Klagenfurt, Austria

**Rohan Byrne**  Transport, Health and Urban Design Research Laboratory, University of Melbourne (UoM), Victoria, Australia

**Ernesto Carella**  University of Oxford, Oxford, UK

**Federico Cecconi**  LABSS (Laboratory of Agent Based Social Simulation), Institute of Cognitive Sciences and Technologies (CNR), Rome, Italy

**Edmund Chattoe-Brown**  School of Media, Communication and Sociology, University of Leicester, Leicester, UK

**Alexis Comber**  School of Geography, University of Leeds, Leeds, UK

**Marco Cremonini**  Department of Social and Political Sciences, University of Milan, Milan, Italy

**Ewa Dabrowska-Prokopowska**  Institute of Sociology, University of Bialystok, Bialystok, Poland

**Carlo Debernardi**  Department of Social and Political Sciences, University of Milan, Milan, Italy

**Mark Dechesne**  Universiteit Leiden, Leiden, The Netherlands

**Mijke van den Hurk**  Utrecht University, Utrecht, The Netherlands; Dutch National Police, The Hague, The Netherlands

**Frank Dignum**  Utrecht University, Utrecht, The Netherlands; Department of Computing Science, Umeå University, Umeå, Sweden; CVUT Prague, Prague, Czech Republic

**Özge Dilaver**  Northumbria University, Newcastle Upon Tyne, UK

**Michael Drexler**  Ocean Conservancy, Washington, D.C., USA

**Bert Droste-Franke**  Institut für qualifizierende Innovationsforschung und -beratung GmbH (IQIB), Bad Neuenahr-Ahrweiler, Germany

**Julie Dugdale**  Grenoble Informatics Laboratory (LIG), Université Grenoble Alpes, Grenoble, France

**Dorine Duives**  Faculty of Civil Engineering and Geosciences, Delft University of Technology, Stevinweg 1, 2628 CN Delft, Netherlands

**Bruce Edmonds**  Centre for Policy Modelling, Manchester Metropolitan University, Manchester, UK

**Maël Franceschetti**  CNRS, Sorbonne Université, Paris, France

**Nobutada Fujii**  Graduate School of System Informatics, Kobe University, Kobe, Hyogo, Japan

**Jiaqi Ge**  School of Geography, University of Leeds, Leeds, UK

**Dario Germani** Department of Political and Social Sciences, University of Cagliari, Cagliari, Italy

**Luca Gerotto** Faculty of Economics, Department of Economics and Finance, Universitá Cattolica del Sacro Cuore, Milan, Italy

**Ross Gore** Virginia Modeling, Analysis and Simulation Center, Old Dominion University, Norfolk, VA, USA

**Arjan Gosal** University of Leeds, Leeds, UK

**Nastasija Grujić** BioSense Institute, University of Novi Sad, Novi Sad, Serbia

**Saida Hachimi El Idrissi** IMI Laboratory, Ibn Zohr University, Agadir, Morocco

**Bezza Hafidi** IMI Laboratory, Ibn Zohr University, Agadir, Morocco

**Lai Kwun Hang** Centre for Science and Technology Studies (CWTS), Leiden University, Leiden, Netherlands

**Etzion Harari** Rafael—Advanced Defense Systems Ltd., Gazit Institute, Tel-Aviv, Israel

**Cédric Herpson** CNRS, Sorbonne Université, Paris, France

**George Hodulik** Center for Mind and Culture, Boston, MA, USA

**Sebastian Hoffmann** Technology Studies Group, Faculty of Social Sciences, TU Dortmund University, Dortmund, Germany

**Sascha Holzhauer** Section Integrated Energy Systems, University of Kassel, Kassel, Germany

**Arvid Horned** Department of Computing Science, Umeå University, Umeå, Sweden

**Jakob Irnich** Faculty of Technology Policy and Management, Delft University of Technology, Jaffalaan 5, 2628 BX Delft, The Netherlands

**Lukas Jansen** Section Integrated Energy Systems, University of Kassel, Kassel, Germany

**Norman L. Johnson** Referentia Systems, Honolulu, HI, USA

**Georg Jäger** University of Graz, Graz, Austria;
Institute of Environmental Systems Sciences, Graz, Austria

**Toshiya Kaihara** Graduate School of System Informatics, Kobe University, Kobe, Hyogo, Japan

**František Kalvas** University of West Bohemia, Plzeň, Czech Republic

**Christian Kammler** Department of Computing Science, Umeå University, Umeå, Sweden

**Jean-Daniel Kant**  CNRS, Sorbonne Université, Paris, France

**Ravshanbek Khodzhimatov**  Digital Age Research Center, University of Klagenfurt, Klagenfurt, Austria

**Marie Lisa Kogler**  Institute of Environmental Systems Sciences, University of Graz, Graz, Austria

**Daisuke Kokuryo**  Graduate School of System Informatics, Kobe University, Kobe, Hyogo, Japan

**Friedrich Krebs**  Section Integrated Energy Systems, University of Kassel, Kassel, Germany;
Fraunhofer Institute for Energy Economics and Energy System Technology, Kassel, Germany

**Piotr Paweł Laskowski**  Institute of Sociology, University of Bialystok, Białystok, Poland

**Sebastian Lehnhoff**  Department of Computing Science, University of Oldenburg, Oldenburg, Germany

**Stephan Leitner**  Department of Management Control and Strategic Management, University of Klagenfurt, Klagenfurt, Austria

**Chunhui Li**  University of Leeds, Leeds, UK

**Philippe Madiès**  Grenoble INP, Université Grenoble Alpes, CERAG, Grenoble, France

**Roderick McClure**  Faculty of Health and Medicine, University of New England, Armidale, NSW, Australia

**Kai G. Mertens**  Institute of Management Accounting and Simulation, Hamburg University of Technology, Hamburg, Germany

**Matthias Meyer**  Institute of Management Accounting and Simulation, Hamburg University of Technology, Hamburg, Germany

**Ruth Meyer**  Centre for Policy Modelling, Manchester Metropolitan University, Manchester, UK

**Yuki Misu**  Waseda University, Tokyo, Japan

**Birgit Müller**  Helmholtz Centre for Environmental Research—UFZ, Leipzig, Germany

**Jean-Pierre Müller**  CIRAD—UMR SENS, Montpellier, France

**Johanna Myrzik**  Institute of Automation, University of Bremen, Bremen, Germany

**Mohamed Nemiche**  Polydisciplinary Faculty of Taza, Taza, Morocco

**Khoa Nguyen**  HES-SO Valais Wallis, SILab, Sierre, Switzerland

**Alain Nkusi**  Teesside University, Middlesbrough, UK

**Joshua Omoju**  Northumbria University, Newcastle Upon Tyne, UK

**Nicolas Payette**  University of Oxford, Oxford, UK

**Paolo Pellizzari**  Department of Economics, Ca' Foscari University of Venice, Venice, Italy

**Valentino Piana**  HES-SO Valais Wallis, SILab, Sierre, Switzerland

**Jannick Plähn**  Institute of Management Accounting and Simulation, Hamburg University of Technology, Hamburg, Germany

**Brian Powers**  Arizona State University, Tempe, USA

**Ivan Puga-Gonzalez**  Center for Modeling Social Systems, NORCE, Kristiansand, Norway

**Hasina Lalaina Rakotonirainy**  Informatique, Géomatique, Mathématiques et Décisions (IGMA), Andrainjato, University of Fianarantsoa, Fianarantsoa, Madagascar

**Tokimahery Ramarozaka**  Informatique, Géomatique, Mathématiques et Décisions (IGMA), Andrainjato, University of Fianarantsoa, Fianarantsoa, Madagascar

**Ashwin Ramaswamy**  Mumbai, India

**Alexandra Rausch**  University of Klagenfurt, Klagenfurt, Austria

**Paul Reaidy**  Grenoble INP, Université Grenoble Alpes, CERAG, Grenoble, France

**Christian Rehtanz**  Institute of Energy Systems, Energy Efficiency and Energy Economics, TU Dortmund University, Dortmund, Germany

**Daniel Reisinger**  Institute of Environmental Systems Sciences, University of Graz, Graz, Austria

**Francesco Renzini**  Department of Social and Political Sciences, University of Milan, Milan, Italy

**Tomer Rokita**  Rafael—Advanced Defense Systems Ltd., Gazit Institute, Tel-Aviv, Israel

**Mariusz Rybnik**  Institute of Computer Science, University of Bialystok, Bialystok, Poland

**Debopama Sen Sarma**  Institute of Energy Systems, Energy Efficiency and Energy Economics, TU Dortmund University, Dortmund, Germany

**Steven Saul**  Arizona State University, Tempe, USA

**Thomas Schmickl**  University of Graz, Graz, Austria;
Institute of Biology, Graz, Austria

**René Schumann**  HES-SO Valais Wallis, SILab, Sierre, Switzerland

**Jan Sören Schwarz**  Department of Computing Science, University of Oldenburg, Oldenburg, Germany

**Davide Secchi**  University of Southern Denmark, Slagelse, Denmark

**Sachith Seneviratne**  Transport, Health and Urban Design Research Laboratory, University of Melbourne (UoM), Victoria, Australia

**Nikitas M. Sgouros**  Department of Digital Systems, University of Piraeus, Piraeus, Greece

**Gaku Shimizu**  Graduate School of System Informatics, Kobe University, Kobe, Hyogo, Japan

**F. LeRon Shults**  Center for Modeling Social Systems, NORCE, Kristiansand, Norway;
Institute for Global Development and Social Planning, University of Agder, Kristiansand, Norway

**Michael D. Slater**  The Ohio State University, Columbus, OH, USA

**Barbara Sonzogni**  Department of Communication and Social Research (CORIS), Sapienza University of Rome, Rome, Italy

**Flaminio Squazzoni**  Department of Social and Political Sciences, University of Milan, Milan, Italy

**John C. Stevenson**  Long Beach Institute, Long Beach, NY, USA

**Emilio Sulis**  Computer Science Department, University of Turin, Turin, Italy

**Timo Szczepanska**  UiT The Arctic University of Norway, Tromsö, Norway

**Shingo Takahashi**  Waseda University, Tokyo, Japan

**Konrad Talmont-Kaminski**  Society and Cognition Unit, University of Bialystok, Bialystok, Poland;
Institute of Sociology, University of Bialystok, Białystok, Poland

**Annina Thaller**  Institute of Environmental Systems Sciences, University of Graz, Graz, Austria

**Jason Thompson**  Transport, Health and Urban Design Research Laboratory, University of Melbourne (UoM), Victoria, Australia

**Eva Margretha Timmer**  Department of Social Sciences, Information Technology Group, Wageningen University and Research, Wageningen, The Netherlands

**Marco Tolotti**  Department of Management, Ca' Foscari University of Venice, Venice, Italy

**Elpida Tzafestas**  Laboratory of Cognitive Science, Department of History and Philosophy of Science, National and Kapodistrian University of Athens, Ano Ilisia, Athens, Greece

**Chrisja Naomi van de Kieft**  Department of Social Sciences, Information Technology Group, Wageningen University and Research, Wageningen, The Netherlands

**H. Van Dyke Parunak**  Parallax Advanced Research, Beavercreek, OH, USA

**Loïs Vanhée**  Department of Computing Science, Umeå University, Umeå, Sweden

**Harko Verhagen**  Department of Computer and Systems Sciences, Stockholm University, Kista, Sweden

**Ben Vermeulen**  Institut für qualifizierende Innovationsforschung und -beratung GmbH (IQIB), Bad Neuenahr-Ahrweiler, Germany

**Katyana Vert-Pre**  Arizona State University, Tempe, USA

**Rajith Vidanaarachchi**  Transport, Health and Urban Design Research Laboratory, University of Melbourne (UoM), Victoria, Australia;
Faculty of Health and Medicine, University of New England, Armidale, NSW, Australia

**Michael Vogrin**  University of Graz, Graz, Austria;
Institute of Biology, Graz, Austria

**Natalie van der Wal**  Faculty of Technology Policy and Management, Delft University of Technology, Jaffalaan 5, 2628 BX Delft, The Netherlands

**Friederike Wall**  Department of Management Control and Strategic Management, University of Klagenfurt, Klagenfurt, Austria

**Yiyu Wang**  School of Geography, University of Leeds, Leeds, UK

**Roseline Wanjiru**  Northumbria University, Newcastle Upon Tyne, UK

**Tom Warendorf**  Institute of Automation, University of Bremen, Bremen, Germany

**Johannes Weyer**  Technology Studies Group, Faculty of Social Sciences, TU Dortmund University, Dortmund, Germany

**Nanda Wijermans**  Stockholm Resilience Centre, Stockholm University, Stockholm, Sweden

**Wesley J. Wildman**  Center for Mind and Culture, Boston, MA, USA;
Boston University, Boston, MA, USA

**Meike Will**  Helmholtz Centre for Environmental Research—UFZ, Leipzig, Germany

**Guilherme Wood**  University of Graz, Graz, Austria;
Institute of Psychology, Graz, Austria

**Marcin Wozniak**  Faculty of Human Geography and Planning, Adam Mickiewicz
University, Poznan, Poland

**Haifeng Zhao**  Transport, Health and Urban Design Research Laboratory, University of Melbourne (UoM), Victoria, Australia

**Guy Ziv**  University of Leeds, Leeds, UK

# The New Frontiers of Social Simulation in the Data Science Era: An Introduction to the Proceedings

**Francesco Renzini, Carlo Debernardi, Federico Bianchi, Marco Cremonini, and Flaminio Squazzoni**

**Abstract** This chapter introduces the proceedings of the Social Simulation Conference 2022 by providing a brief overview of the impact of social simulation in various research areas. By focusing on the key role of agent-based modeling, we argue that social simulation has a unique position in the wider data science area. This is because it can enrich the predominantly inductive, data-driven, pattern oriented approach of computational social science with deductive, hypothesis-driven, explanatory, mechanism-detection models. Furthermore, social simulation can also work in areas and for contexts where data is not available, experiments cannot be performed or in which scenario exploration is paramount. We would also like to focus on areas and aspects where methodological improvement and cross-methodological integration are required to enhance the potential of social simulation in various communities. In the final section, we introduce the structure and sections of the proceedings.

F. Renzini · C. Debernardi · F. Bianchi · M. Cremonini · F. Squazzoni (✉)
Department of Social and Political Sciences, University of Milan, via Conservatorio 7, 20122, Milan, Italy
e-mail: flaminio.squazzoni@unimi.it

F. Renzini
e-mail: francesco.renzini@unimi.it

C. Debernardi
e-mail: carlo.debernardi1@unimi.it

F. Bianchi
e-mail: federico.bianchi1@unimi.it

M. Cremonini
e-mail: marco.cremonini@unimi.it

# 1   The Computational Social Science Landscape

Recently, there has been increased academic recognition of social simulation in various research areas indicating a relatively coherent interdisciplinary field, where similar modeling frameworks and tools are used to study complex social dynamics and processes [64]. This is testified by the number of reviews on the use of agent-based models recently published in Economics [3], Management [71], and Epidemiology and Public Health [37, 68]. There has been a specific focus also on the recent COVID-19 pandemic [43], and reviews in Ecology [33], Political Science [13], Sociology [25], Criminology [34], as well as in fields such as opinion dynamics and social influence [26], social identity [59], peer review [25], land use [47], urban residential patterns [36], and socio-environmental systems [19].

Besides testifying to the general interest towards social simulation models, these overviews on the state-of-the-art have identified certain common strengths of this type of research in response to various domain-specific traditions, priorities, and challenges. First, by considering agent and rule heterogeneity, agent-based models have allowed us to explore the underlying mechanisms of non-linear, emergent, and complex dynamics overcoming certain empirically implausible assumptions required by conventional mathematical models to ensure analytical tractability [24, 31, 61]. For instance, in Economics and Finance, relaxing the typical assumptions of market equilibrium, representative agents and rational expectations was pivotal to examine certain empirical puzzles not appropriately explained by conventional models, such as business cycle fluctuations [14], bankruptcy avalanches [15], and the underlying mechanisms behind financial crises [7], including the role of housing markets as determinants of financial instability and contagion [28].

In order to show the importance of heterogeneity also for macroeconomic forecasting, Poledna and Miess developed a calibrated computational model with a population of a million heterogeneous agents [50]. Their results have competed with standard modelling techniques employed by economic institutions. Furthermore, by considering heterogeneously mixed populations, epidemiologists and public health researchers also improved their understanding of the unfolding of outbreaks and tested interventions to improve collective outcomes in artificial scenarios [37, 68]. They were also able to examine the global effect of neighborhood-level patterns of healthy behavior [1].

Secondly, social simulation models have allowed us to consider the effects of non-trivial interaction structures such as complex spatial or social network structures, on aggregate dynamics due to individual adaptive behavior. Recently, this has allowed epidemiologists to gain leverage on individual mobility and social network data to forecast outbreak dynamics [27]. For instance, Eubank et al. studied the diffusion of an infectious disease on spatially bounded social networks using social and geographic contact networks, calibrated from urban transportation grids of various US cities, census data, and surveys on mobility patterns [23]. Manzo and van de Rijt calibrated an agent-based model with heterogeneous network structures from survey

data to assess the effect of interventions targeted to network hubs to reduce the spread of COVID-19 ([45] see also [69]).

Thanks to their dynamic nature [13], social simulation models have also permitted researchers to study the aggregate effect of causal feedback among populations of adapting agents. For instance, Epstein et al. [22] suggested that epidemiological models did not succeed in estimating the diffusion of fatal diseases, such as AIDS, because they followed too simplistic assumptions on social behavior, including perfect mixing [37]. This explains why computational models with adaptive agents are now popular in epidemiology and public health [46].

At the intersection of social network analysis and social simulation, researchers have begun to use agent-based models to examine the role of generative behavioral mechanisms on social network formation, such as advice and friendship networks, by exploring various underlying mechanisms and fitting models to empirical data [60, 66]. In more qualitative research areas, following the famous application on the Anasazi in the Long House Valley in Arizona between 800 and 1350 [2], archaeologists are now using rich, context-specific agent-based models to test the effect of social structures on past trajectories of extinct populations in various geographical regions (e.g., see [56]). Re-running past historical periods with computer simulations requires integrating data and insights from various disciplines so that counter factual 'as if' tests on historical trajectories can be performed [9].

Social simulation models are also useful when data is not available, i.e., when collecting data or running experiments is unethical, unfeasible or impractical. Examples can be found again in Epidemiology [68], but also in Criminology [34], Science of Science research on the effect of peer review on publication bias (e.g., [5, 25]), and Management [71]. Similar models have also been used to test policies in artificial scenarios before any direct intervention—often too costly or unfeasible due to ethical constraints. For instance, in Macroeconomics, scholars tested the impact of fiscal [48], monetary [57] and labor policies [12], as well as the coupled climate and economic assessment of complex, long-term outcomes [39], while considering heterogeneous responses and non-linear interactions. Similar ex-ante policy assessment with agent-based models has been explored in the field of refugee crisis management [67].

While outlining either general or context-specific strengths of social simulation models, these overviews have also provided a systematic assessment of certain perceived weaknesses or areas of necessary improvement. First, the possibility of modeling agent heterogeneity is a double-edged sword. For example, abandoning abstract, simple and generally accepted theories for micro or meso realism often leads modelers to theoretical bricolage or excessive freedom on important model building blocks (e.g., see [55]). It also exacerbates the challenge of parameter calibration and the risk of overfitting [3, 13, 49, 50].

Excessive model idiosyncrasies undermine the establishment of common model building protocols, thus compromising comparison and reproducibility (e.g., see [26]). The high exploration cost of large parameter spaces (the curse of dimensionality; [3]) increases model dependence on data availability for parameter calibration. However, even if sufficient data is available to calibrate models with many

parameters, there is a problem of *appropriateness* in that calibrated data could come from observational studies with context-specific underlying causal mechanisms [68]. Furthermore, even with sufficient and appropriate data, calibration is never "a free lunch" [8], in that the choice of estimation algorithms or summary statistics for an empirical target can greatly depend on relatively inscrutable heuristics judgement. Finally, empirical calibration is computationally demanding also due to the inherent stochasticity of agent-based models [34, 68]. While High-Performance Computer (HPC) clusters are increasingly available [53], computational social scientists are often disadvantaged in their access compared to other research communities [51].

Another perceived weakness is that empirically calibrated models tend to produce high volumes of output data, thus making the estimation of the effect of highly non-linear, path-dependent social dynamics even more critical [41]. In principle, there is not even any guarantee that well-calibrated and validated models would satisfy sufficient standards of out-of-sample validity [71]. Insufficient model documentation and lack of transparency are also perceived as additional problems affecting the credibility of simulation model outputs, even if significant steps forward have been made regarding replicability and transparency, e.g., with the ODD (Overview, Design concepts, and Details) protocol [32]. Finally, in certain reviews, the under-development of equity considerations in social simulation, the lack of clarity on model purposes, and the rare acknowledgement of positionality and implicit bias of the modeler, are all seen as social simulation areas in which improvement is necessary [18, 72].

## 2 The Role of Agent-Based Modeling in the Data Science Era

In the 1990s, prior to the advent of big data, the label "computational social sciences" mostly meant the application of computer simulation models in the social sciences, with a central role played by agent-based models [10]. Computational models were especially used to simulate complex social systems from first behavioral principles (e.g., see [21]). Theory was dominant over empirical data [16, 58], and computational modeling was mostly considered as a data-generating tool and a method to test hypotheses on social dynamics via artificial simulation experiments [62].

Now, the situation has changed as does the perceived meaning of computational social sciences. First, there has been an empirical turn in computational modeling fueled by the advent of big data and machine learning [40]. Models are now mostly used to detect hidden patterns in large-scale data, enhance inference from data and facilitate prediction. The emphasis is on the virtue of data, with computation now mostly used for data-driven explorations [17]. Second, even when theory is used to inform computational models, the dominant framework involves concepts from physics that emphasize structural properties and aggregate distributions (e.g., power laws). Insights from social and behavioral sciences that point to identify underlying generative mechanisms are less relevant [6]. For instance, one of the recent Horizon EU calls for proposals included a footnote reporting that computational social

sciences were to be intended as "methods developed in statistical physics to take advantage of the very rich big data sets".[1] This was not the definition proposed in a manifesto for computational social sciences previously published by a team of agent-based computational modelers, where theory-driven behavioral models were still seen as central for computational social sciences [11].

However, while stimulating research in this area, this new trend of data-driven computational social science has left certain unsolved questions, which, in our opinion, can be considered at the core of social simulation models. First, although important, drawing statistical inferences from large-scale data does not necessarily mean being able to identify causal generative mechanisms behind emergent social patterns, especially in cases where individual preferences and opportunities cannot be observed empirically or tested experimentally. As suggested by previous research [29, 38, 44, 68], agent-based modeling is suitable for exploring multiple causal paths and addressing problems of multiple realizability (i.e., the same effect generated by different social causes and paths), thanks to systematic experiments on the effect of varying initial conditions and parameter values on observed outcomes. This is key to exploring potentially alternative generative mechanisms of the same empirically detected pattern also via counterfactual testing [63].

It is also worth noting that some scholars in artificial intelligence and machine learning have started to reflect on the key role of causal analysis and the danger of black-box explanations (e.g., see [42, 54]). However, in order to play a role here, the complexity of social simulation models in terms of empirically calibrated population size and attributes, time and spatial scales of interaction and statistical treatment of stochastic elements of any social complex system must be taken to the next level.

Second, while detecting aggregate patterns from large-scale data is key both for testing existing theories and measurements and contributing to generate new ones, research on social systems' behavior also requires scenario analysis and experimentation, which are difficult in data-driven computational social science [30]. There is a strategic intersection between policy, experimentation and models for which the position of computational social sciences greatly depends on behavioral, experimental and computer simulation research [63]. Lack of transparency in data science research, ethical implications of large-scale behavioral policy experiments, and the request for context-specific 'what if' scenario analysis are all factors enhancing the key role of empirically-calibrated, context-specific, social simulation models in the computational social science toolkit [65].

This brief overview on the strengths and weaknesses of social simulation models suggests that there is a need to integrate various methods and approaches and increase the connection between data and theory through computational social science models. Examples of this integration exist in research areas on collective action and the management of common-pool resources, where agent-based models have been used to integrate findings from behavioral experiments and context-tracing qualitative

---

[1] This text refers to the call "Past, present and future of democracies" (HORIZON-CL2-2024-DEMOCRACY-01) and can be found at this link: https://ec.europa.eu/info/funding-tenders/opportunities/portal/screen/opportunities/topic-details/horizon-cl2-2024-democracy-01-06.

analysis [52]. This cross-methodological practice should also be applied within the broad field of computational social science.

For instance, in a recent special section of *JASSS-Journal of Artificial Societies and Social Simulation* on inverse generative social sciences (*iGSS*), Epstein has recommended the application of genetic algorithms to discover the specific generative behavioral rules that would lead a population of agents to reproduce empirically observed target distributions [20]. Brilliant examples of these applications included the study of the effect of norms on drinking behaviour [70] and behavioral rules that determine mixed-residential segregation patterns [35]. Given that the idea of *iGSS* leverages on genetic search algorithms to find potential explanations for large-scale data patterns, modelers can concentrate on the agent-level constituents of such explanations and their permissible combinations. Not only does this integration of behavioral constituents, empirical data and artificial intelligence lead to the development of new theories, which even the brightest modeler could not develop simply by intuition; this integration could also guide future data collection or offer new ways to detect patterns in existing large-scale data.

In conclusion, the overview on the state-of-the-art of social simulation and a brief analysis on its role in the current landscape of computational social science suggest a tighter cross-fertilization between hypothesis-driven models and data collection, which aim to identify causal generative mechanisms and perform counter factual, data-driven analysis which can identify explanatory patterns in fine-grained, large-scale data. In our opinion, the social simulation community can support this integration if it continues to be an open, cross-disciplinary, trans-domain hub and be capable of refining its specificity in an unpredictable world where research, data and methods are rapidly changing.

## 3   The Book Content

The remainder of this book includes the post-proceedings of the Social Simulation Conference 2022 and is organized into five thematic sections. Each section offers contributions that reflect the multifaceted nature of social simulation as a cross-disciplinary research area.

Section 1—Social Behavior—includes contributions that consider the challenge of understanding social behavior in various research areas. Examples range from the dynamics of societies' formation to social learning, social norms and decision-making in strategic social contexts. Some chapters are more applied and public decision making-oriented, including pedestrian behavior and evacuation dynamics, the diffusion of infectious disease diffusion and the endogenous relation of fisheries and fish reproduction.

Section 2—Social Identity and Social Influence—includes chapters that focus on the importance of social networks in shaping individual identity. They provide a link between social simulation, social psychology and network studies. Examples range from polarization and political radicalization to friendship formation.

Section 3—Management and Economics—includes contributions that explore the application of agent-based modeling to various relevant issues in management and economics. Examples include topics such as: process mining, task allocation, replication control, honesty norms, innovation systems, wine pricing policies, credit relationship dynamics, decision-making, virtual teams and the equity premium puzzle. These contributions highlight the versatility and usefulness of agent-based modeling to examine complex dynamics often difficult to consider within standard, disciplinary modeling frameworks.

Section 4—Energy and Climate— explores the application of agent-based modeling to examine sustainability and energy-related issues. The chapters cover relevant topics such as the adoption of agri-environment schemes by farmers, the co-simulation of socio-technical energy systems, investments in heating technology, the impact of beliefs on food and climate change on dietary adoption and public acceptance of green mobility policies. These contributions testify to the growing importance of social simulation in research areas where estimating large-scale, aggregate outcomes in complex contexts is paramount as behaviors, interactions and environments play a crucial role in emergent, non-linear dynamics.

Finally, Sect. 5—Tools and Methods—includes contributions that provide useful insights to the community of social simulation scholars by discussing the process of model development both theoretically and practically. The chapters explore relevant topics, such as the impact of modeling choices and the link between theory and data, and provide examples of important modeling techniques.

## References

1. Auchincloss, A.H., Diez Roux, A.V.: A new tool for epidemiology: the usefulness of dynamic-agent models in understanding place effects on health. Am. J. Epidemiol., 16811–16818 (2008)
2. Axtell, R.L., Epstein, J.M., Dean, J.S., Gumerman, G.J., Swedlund, A.C., Harburger, J., Parker, M.: Population growth and collapse in a multi-agent model of the Kayenta Anasazi in Long House Valley. Proc. Nat. Acad. Sci. **99**(3), 7275–7279 (2002)
3. Axtell, R.L., Farmer, J.D.: Agent-based modeling in economics and finance: past, present, and future. J. Econ. Liter. Forthcoming (2023)
4. Bianchi, F., Squazzoni, F.: Agent-based models in sociology. Wiley Interdiscip. Rev. Comput. Statist. **7**(4), 284–306 (2015)
5. Bianchi, F., Squazzoni, F.: Can transparency undermine peer review? A simulation model of scientist behavior under open peer review. Sci. Publ. Policy **49**(5), 791–800 (2022)
6. Boero, R.: Behavioral Computational Social Science. Hoboken, NJ Wiley (2015)
7. Bookstaber, R.: Agent-based models for financial crises. Ann. Rev. Fin. Econ. **9**, 85–100 (2017)
8. Carrella, E.: No free lunch when estimating simulation parameters. J. Artif. Soc. Social Simul. **24**(2), 7 (2021). Retrieved from: https://jasss.soc.surrey.ac.uk/24/2/7.html
9. Chliaoutakis, A., & Chalkiadakis, G.: An agent-based model for simulating inter-settlement trade in past societies. J. Artif. Soc. Soc. Simul. **23**(3), 10 (2020). Retrieved from http://jasss.soc.surrey.ac.uk/23/3/10.html
10. Cioffi-Revilla, C.: Introduction to Computational Social Science. Principles and Applications. Springer Verlag, Berlin, Heidelberg (2014)
11. Conte, R., Gilbert, N., Bonelli, G., Cioffi-Revilla, C., Deffuant, G., Kertesz, J., et al.: Manifesto of computational social science. Euro. Phys. J. Special Top. **214**, 325–346 (2012)

12. Dawid, H., Gemkow, S., Harting, P., Neugart, M.: Labor market integration policies and the convergence of regions: the role of skills and technology diffusion. J. Evol. Econ. **22**, 543–562 (2012)
13. De Marchi, S., Page, S.E.: Agent-based models. Ann. Rev. Polit. Sci. **17**, 1–20 (2014)
14. Delli Gatti, D., Di Guilmi, C., Gaffeo, E., Giulioni, G., Gallegati, M., Palestrini, A.: A new approach to business fluctuations: heterogeneous interacting agents, scaling laws and financial fragility. J. Econ. Behav. Org. **56**(4), 489–512 (2005)
15. Delli Gatti, D., Gallegati, M., Greenwald, B.C., Russo, A., Stiglitz, J.E.: Business fluctuations and bankruptcy avalanches in an evolving network economy. J. Econ. Interaction Coordination **4**, 195–212 (2009)
16. De Marchi, S.: Computational and Mathematical Modeling in the Social Sciences. Cambridge University Press (2005)
17. Edelmann, A.,Wolff, T., Montagne, D., Bail, C.A.: Computational social science and sociology. Ann. Rev. Sociol. **46**(1), 61–81 (2020)
18. Edmonds, B., Le Page, C., Bithell, M., Chattoe-Brown, E., Grimm, V., Meyer, R., Squazzoni, F.: Different modelling purposes. J. Artif. Soc. Soc. Simul. **22**(3), 6 (2019). Retrieved from http://jasss.soc.surrey.ac.uk/22/3/6.html
19. Elsawah, S., Filatova, T., Jakeman, A.J., Kettner, A.J., Zellner, M.L., Athanasiadis, I.N., Lade, S.J.: Eight grand challenges in socio-environmental systems modeling. Socio-Environ. Syst. Modell. **2**, 16226 (2020)
20. Epstein, J.M.: Inverse generative social science: backward to the future. J. Artif. Soc. Soc. Simul. **26**(2), 9 (2023). Retrieved from http://jasss.soc.surrey.ac.uk/26/2/9.html
21. Epstein, J.M., Axtell, R.: Growing Artificial Societies. Social Science from the Bottom Up. The MIT Press, Cambridge, MA (1996)
22. Epstein, J.M., Parker, J., Cummings, D., Hammond, R.A.: Coupled contagion dynamics of fear and disease: mathematical and computational explorations. PLoS One **3**(12), 1–11 (2008)
23. Eubank, S., Guclu, H., Anil Kumar, V.S.: Modelling disease outbreaks in realistic urban social networks. Nature **429**, 180–184 (2008)
24. Fagiolo, G., Roventini, A.: Macroeconomic policy in DSGE and agent-based models redux: new developments and challenges ahead. J. Artif. Soc. Soc. Simul. **20**(1), 1 (2017). Retrieved from https://www.jasss.org/20/1/1.html
25. Feliciani, T., Luo, J., Ma, L., Lucas, P., Squazzoni, F., Marušić, A., Shankar, K.: A scoping review of simulation models of peer review. Scientometrics **121**, 555–594 (2019)
26. Flache, A., Mäs, M., Feliciani, T., Chattoe-Brown, E., Deffuant, G., Huet, S., Lorenz, J.: Models of social influence: towards the next frontiers. J. Artif. Soc. Soc. Simul. **20**(4), 2 (2017)
27. Friás-Martínez, E., Williamson, G., Friás-Martínez, V.: An agent-based model of epidemic spread using human mobility and social network information. In: IEEE Conference on Social Computing (2011)
28. Geanakoplos, J., Axtell, R., Farmer, J.D., Howitt, P., Conlee, B., Goldstein, J., Yang, C.Y.: Getting at systemic risk via an agent-based model of the housing market. Am. Econ. Rev. **102**(3), 53–58 (2012)
29. Giabbanelli, P.J., Tison, B., Keith, J.: The application of modeling and simulation to public health: assessing the quality of agent-based models for obesity. Simul. Modell. Pract. Theor. **108**, 102268 (2021)
30. Gilbert, N., Ahrweiler, P., Barbrook-Johnson, P., Narasimhan, K.P., Wilkinson, H.: Computational modelling of public policy: reflections on practice. J. Artif. Soc. Soc. Simul. **21**(1), 14 (2018). Retrieved from http://jasss.soc.surrey.ac.uk/21/1/14.html
31. Gilbert, N., Troitzsch, K.G.: Simulation for the Social Scientist. Open University Press, Milton Keynes (2005)
32. Grimm, V., Berger, U., DeAngelis, D.L., Polhill, J.G., Giske, J., Railsback, S.F.: The ODD protocol: a review and first update. Ecol. Modell. **221**(23), 2760–2768 (2010)
33. Grimm, V., Railsback, S.F.: Agent-based models in ecology: patterns and alternative theories of adaptive behaviour. In: Billari, F.C., Fent, T., Prskawetz, A., Scheffran, J. (eds.) Agent-Based Computational Modelling. Springer, Berlin/Heidelberg (2006)

34. Groff, E.R., Johnson, S.D., Thornton, A.: State of the art in agent-based modeling of urban crime: an overview. J. Quantitat. Criminol. **35**, 155–193 (2019)
35. Gunaratne, C., Hatna, E., Epstein, J.M., Garibay, I.: Generating mixed patterns of residential segregation: an evolutionary approach. J. Artif. Soc. Soc. Simul. **26**(2), 7 (2023). Retrieved from https://www.jasss.org/26/2/7.html
36. Huang, Q., Parker, D.C., Filatova, T., Sun, S.: A review of urban residential choice models using agent-based modeling. Environ. Plann. B Plann. Des. **41**, 661–689 (2014)
37. Hunter, E., MacNamee, B., Kelleher, J.D.: A taxonomy for agent-based models in human infectious disease epidemiology. J. Artif. Soc. Soc. Simul. **20**(3), 2 (2017). Retrieved from https://www.jasss.org/20/3/2.html
38. Klein, M., Frey, U.J., Reeg, M.: Models within models—agent-based modelling and simulation in energy systems analysis. J. Artif. Soc. Soc. Simul. **22**(4), 6 (2019). Retrieved from http://jasss.soc.surrey.ac.uk/22/4/6.html
39. Lamperti, F., Dosi, G., Napoletano, M., Roventini, A., Sapio, A.: Faraway, so close: coupled climate and economic dynamics in an agent-based integrated assessment model. Ecol. Econ. **150**, 315–339 (2018)
40. Lazer, D., Pentland, A., Adamic, L., Aral, S., Barabási, A.-L., Brewer, D., et al.: Computational social science. Science **323**(5915), 721–723 (2009)
41. Lee, J.S., Filatova, T., Ligmann-Zielinska, A., Hassani-Mahmooei, B., Stonedahl, F., Lorscheid, I., Parker, D.C.: The complexities of agent-based modeling output analysis. J. Artif. Soc. Soc. Simul. **18**(4), 4 (2015). Retrieved from https://www.jasss.org/18/4/4.html
42. Lipton, Z.C.: The mythos of model interpretability: in machine learning, the concept of interpretability is both important and slippery. Queue **16**(3), 31–57 (2018, June)
43. Lorig, F., Johansson, E., Davidsson, P.: Agent-based social simulation of the Covid-19 pandemic: a systematic review. J. Artif. Soc. Soc. Simul. **24**(3), 5 (2012). Retrieved from https://www.jasss.org/24/3/5.html
44. Manzo, G.: Agent-Based Models and Causal Inference. Hoboken, NJ Wiley (2022)
45. Manzo, G., van de Rijt, A.: Halting SARS-COV-2 by targeting high-contact individuals. J. Artif. Soc. Soc. Simul. **23**(4), 10 (2020). Retrieved from http://jasss.soc.surrey.ac.uk/23/4/10.html
46. Marshall, B.D.: Formalizing the role of agent-based modeling in causal inference and epidemiology. Am. J. Epidemiol. **181**(2), 92–99 (2015)
47. Matthews, R.B., Gilbert, N.G., Roach, A., Polhill, J.G., Gotts, N.M.: Agent-based land-use models: a review of applications. Landsc. Ecol. **22**, 1447–1459 (2007)
48. Napoletano, M., Gaffard, J.-L., Roventini, A.: Time-varying fiscal multipliers in an agent-based model with credit rationing. Working Paper Series 2015/19, Laboratory of Economics and Management (LEM), Scuola Superiore Sant'Anna, Pisa, Italy (2015)
49. Platt, D.: A comparison of economic agent-based model calibration methods. J. Econ. Dyn. Control **113**, 103859 (2020)
50. Poledna, S., Miess, M.G.: Economic forecasting with an agent-based model. Euro. Econ. Rev. **151**, 104306 (2023)
51. Polhill, J.G.: Antisocial simulation: using shared high-performance computing clusters to run agent-based models (2022). Available at https://rofasss.org/tag/garypolhill/
52. Poteete, A.R., Janssen, M.A., Ostrom, E.: Working Together: Collective Action, The Commons, and Multiple Methods in Practice. Princeton, NJ Princeton University Press (2010)
53. Reuillon, R., Leclaire, M., Rey-Coyrehourcq, S.: OpenMOLE, a workflow engine specifically tailored for the distributed exploration of simulation models. Futur. Gener. Comput. Syst. **29**(8), 1981–1990 (2015)
54. Ribeiro, M.T., Singh, S., Guestrin, C.: Why Should I Trust You? Explaining the Predictions of any Classifier (2016). Available at https://arxiv.org/pdf/1602.04938.pdf
55. Richiardi, M., Leombruni, R., Saam, N., Sonnessa, M.: A common protocol for agent-based social simulation. J. Artif. Soc. Soc. Simul. **9**(1), 15 (2006). Retrieved from https://www.jasss.org/9/1/15.html

56. Romanowska, I., Wren, C.D., Crabtree, S.A.: Agent-Based Modeling for Archaeology: Simulating the Complexity of Societies. Santa Fe Institute Press (2021)
57. Salle, I., Yıldızoğlu, M., Senegas, M.-A.: Inflation targeting in a learning economy: an ABM perspective. Econ. Modell. **34**, 114–128 (2013)
58. Sawyer, R.K.: Social Emergence: Societies as Complex Systems Cambridge: University Press (2005)
59. Scholz, G., Wijermans, N., Paolillo, R., Neumann, M., Masson, T., Chappin, E., Kocheril, G.: Social agents? A systematic review of social identity formalizations. J. Artif. Soc. Soc. Simul. **26**(2), 56 (2023). Retrieved from https://www.jasss.org/26/2/6.html
60. Snijders, T.A.B., Steglich, C.E.G.: Representing micro-macro linkages by actor-based dynamic network models. Sociol. Methods Res. **44**(2), 222–271 (2015)
61. Squazzoni, F.: The impact of agent-based models in the social sciences after 15 years of incursion. History Econ. Ideas **18**(2), 197–233 (2010)
62. Squazzoni, F.: Agent-Based Computational Sociology. Wiley, Hoboken, NJ (2012)
63. Squazzoni, F., Bianchi, F.: Exploring interventions on social outcomes with in silico, agent-based experiments. In: Damonte, A., Negri, F. (eds.) Causality in policy studies: a pluralist toolbox, pp. 217–234. Springer International Publishing, Cham (2023)
64. Squazzoni, F., Jager, W., Edmonds, B.: Social simulation in the social sciences: a brief overview. Soc. Sci. Comput. Rev. **32**(3), 279–294 (2014)
65. Squazzoni, F., Polhill, G., Edmonds, B., Ahrweiler, P., Antosz, P., Scholz, G., Gilbert, N.: Computational models that matter during a global pandemic outbreak: a call to action. J. Artif. Soc. Soc. Simul. **23**(2), 10 (2020). Retrieved from http://jasss.soc.surrey.ac.uk/23/2/10.html
66. Steglich, C.E.G., Snijders, T.A.B.: Stochastic network modeling as generative social science. In: Gërxhani, K., de Graaf, N., Raub, W. (eds.) Handbook of Sociological Science, pp. 73–95. Edward Elgar Publishing, Northampton, MA (2022)
67. Suleimenova, D., Groen, D.: How policy decisions affect refugee journeys in South Sudan: a study using automated ensemble simulations. J. Artif. Soc. Soc. Simul. **23**(1), 2 (2020). Retrieved from http://jasss.soc.surrey.ac.uk/23/1/2.html
68. Tracy, M., Cerdá, M., Keyes, K.M.: Agent-based modeling in public health: current applications and future directions. Ann. Rev. Publ. Health **39**, 77–94 (2018)
69. Vermeulen, B., Müller, M., Pyka, A.: Social network metric-based interventions? Experiments with an agent-based model of the COVID-19 pandemic in a metropolitan region. J. Artif. Soc. Soc. Simul. **24**(3), 6 (2021). Retrieved from http://jasss.soc.surrey.ac.uk/24/3/6.html
70. Vu, T.M., Buckley, C., Duro, J.A., Brennan, A., Epstein, J.M., Purshouse, R.C.: Can social norms explain long-term trends in alcohol use? Insights from inverse generative social science. J. Artif. Soc. Soc. Simul. **26**(2), 4 (2023). Retrieved from https://www.jasss.org/26/2/4.html
71. Wall, F.: Agent-based modeling in managerial science: an illustrative survey and study. Rev. Managerial Sci. **10**, 135–193 (2016)
72. Williams, T.G., Brown, D.G., Guikema, S.D., Logan, T.M., Magliocca, N.R.: Integrating equity considerations into agent based modeling: a conceptual framework and practical guidance. J. Artif. Soc. Soc. Simul. **25**(3), 1 (2022). Retrieved from https://www.jasss.org/25/3/1.html

# Social Behaviour

# A Cognitive Model of Epistemic Vigilance in Situations of Varying Competence, Consistency, and Utility

**Piotr Paweł Laskowski** ⬛, **Ivan Puga-Gonzalez** ⬛, **F. LeRon Shults** ⬛, and **Konrad Talmont-Kaminski** ⬛

**Abstract** This paper outlines a computational, cognitive model representing how humans may use epistemic vigilance to evaluate socially-provided information in a way that reacts flexibly to differences in the reliability of content versus source vigilance strategies. Furthermore, the model explores how the system reacts in situations where the utility of the information provided is either unrelated to its accuracy or, even, is inversely proportional to it. We find that even a simple model is able to react flexibly to variation in these parameters, providing a basis for further exploration of the phenomenon.

**Keywords** Computational cognitive model · Cognition · Epistemic vigilance · Human cooperation · Bounded rationality

## 1 Introduction

According to epistemic vigilance theory [1], humans use two main methods in the attempt to ensure that the information they learn from others is trustworthy. The first of these consists in vigilance to the content of the information, i.e. the evaluation of the overall plausibility of the information, which can depend upon such things as its coherence with existing beliefs as well as its internal consistency. At the very least, this requires making a quick judgment as to the prima facie plausibility of the content, which may well be subject to any of a number of cognitive biases. The second concerns vigilance to the source of the information, i.e. whether the person providing the information is themselves plausible, which will depend upon

P. P. Laskowski (✉) · K. Talmont-Kaminski
Institute of Sociology, University of Bialystok, Białystok, Poland
e-mail: pio.laskowski@gmail.com

I. Puga-Gonzalez · F. L. Shults
Center for Modeling Social Systems, Norwegian Research Center, Kristiansand, Norway

considerations of their knowledge ability in the area, as well as their attitude or potential conflicts of interest. At a minimum, source vigilance relies upon the social status of the source, thereby making the assumption that social status is related to ability to provide accurate information.

Each of the methods has its own strengths and limitations and is therefore appropriate for use in particular contexts. It is highly significant, therefore, to what degree people are thought to be able to use these two forms of vigilance flexibly, applying them as needed in the situation they find themselves in. In particular, source vigilance is more appropriate where a person has very little ability to judge the content (low competence) and content vigilance is the preferable approach to take where there is little relationship between the apparent plausibility of the source and the quality of the information they offer (low consistency). Listening to one's doctor may be considered an example of the first kind of situation, while choosing one's life insurance on the basis of the actor who advertises it is an example of the latter. What makes it worse is that it is not always clear cut what kind of an environment one is in, so that it becomes necessary to judge on the basis of the outcomes of one's previous decisions.

The issue is further complicated by the question of the relationship between the accuracy of the information one is presented with and its utility. In general, it is important to act upon accurate information—it is no use arriving at the airport to catch the 11 am flight if it actually was scheduled to leave at 9 am. But the relationship between utility and accuracy of information need not be so straight-forward. In extreme cases, it may be that inaccurate information has a higher utility than accurate information. For example, it may be that false beliefs about potential punishment can motivate a group to cooperate to their mutual benefit [2].

We decided to model such a situation in order to determine how a cognitive model reacts to, on the one hand, variation in the relationship between accuracy and utility and, on the other hand, variation in its competence to judge the content of the information as opposed to the consistency between the plausibility of the sources and the quality of the information they offer [3]. The model represents the cognitive process of an agent that is being successively presented with pieces of information. Based on this cognitive process, the agent may provisionally accept or reject the information on the basis of source or content vigilance. If it provisionally accepts the information, it acts upon it, thereby generating utility that may or may not satisfice it [4]. If it satisfices, it accepts the piece of information and strengthens its preference for the form of vigilance it used to evaluate the information. If it does not satisfice, it rejects the piece of information, retaining its prior belief, and weakens its preference for the form of vigilance it used to evaluate it. Having built this model, we then examine how the agent deals with three environments, as relating to the relationship between accuracy and utility of beliefs: (1) where utility is equal to accuracy, (2) where utility is random and unrelated to accuracy, and (3) where utility has an inverse relationship to accuracy.

## 2   Methods

### 2.1   Model Description

The model represents an agent presented with a series of beliefs which they may choose to act upon, where each belief has a degree of accuracy (a value from a uniform distribution [0,1]). The decision whether to act upon the belief presented is made on the basis of either source or content vigilance.

When the belief is judged on the basis of content vigilance, the agent applies its competence (a value between 0 and 1) in order to arrive at an estimate of the accuracy of the content of the belief according to equation 1.

$$EA = N(0, 1, \mu = \text{Accuracy}, \sigma = 1 - \text{Competence}) \qquad (1)$$

where EA represents the estimated accuracy, a value drawn from a truncated normal distribution at [0,1] with mean equal to the accuracy of the belief and standard deviation equal to one minus competence. Thus, the higher the competence of the individual, the more likely that the EA represents a value close to the real accuracy.

When the belief is judged on the basis of source vigilance, the agent relies upon the consistency, in their environment, of the relationship between source quality and information quality (also a value between 0 and 1). When consistency is high, the accuracy of the belief is closely connected to the plausibility of the source presenting it. Since the model does not include actual sources, EA is deemed to be equal to 'source plausibility' calculated 'backwards' from belief accuracy and the consistency of the environment, as per equation 2.

$$EA = N(0, 1, \mu = \text{Accuracy}, \sigma = 1 - \text{Consistency}) \qquad (2)$$

where EA represents the estimated accuracy, a value drawn from a truncated normal distribution at [0,1] with mean equal to the accuracy of the belief and standard deviation equal to 1 minus the consistency of the environment. Thus, the higher the consistency of the environment the more likely that the EA represents a value close to the real accuracy.

At the beginning of the simulation the probabilities of the individual judging the belief on the basis of the content or source are both 0.5. These values may later change depending on which strategy (content or source) leads to higher utility beliefs.

After calculating EA on the basis of the content or the source, the individual tries the belief if the estimated accuracy is higher than the previous estimated accuracy minus an error; where error is kept constant at 0.1 and is necessary to allow the agent to try beliefs of similar but somewhat lower estimated accuracy. If the individual decides to try the belief, then the agent calculates the probability that the outcome satisfices it, according to equation 3:

$$Prob\_satisfaction = \frac{1}{1 + e^{-\lambda * (Utility - 0.5)}} \qquad (3)$$

This equation represents a sigmoidal curve, where lambda gives the shape of the curve and Utility $-$ 0.5 gives the inflection point of the curve. Hence, if utility is equal to 1, the inflection point of the sigmoidal curve will be exactly at 0.5 on the x and y axes. We keep lambda fixed at 10.

We produce three different scenarios regarding the relationship between the value of utility and accuracy. In the first, utility is equal to accuracy, so that beliefs with high accuracy will also have a higher chance of satisficing if they are tried. In the second scenario, utility value is not tied to accuracy. In that case, utility is drawn from a uniform distribution [0,1] each time. Hence, no matter the level of accuracy of the belief, the satisficing level it produces is decided at random. In the third scenario, we reverse the relationship between accuracy and utility from that in the first scenario. Thus, beliefs with low accuracy will have higher utility and vice versa. This means that a low accuracy belief will have a high probability of satisficing.

Once the probability of satisficing is calculated, it is compared to a random number between [0,1], if it does satisfice (i.e. probability of satisficing is higher than the random number), then the agent accepts the belief and increases the probability of using the strategy (content or source) that produced this outcome. For instance, if the belief was judged on the basis of content, then the probability of using this strategy with future beliefs increases according to equation 4:

$$\text{Prob\_content}_{t+1} = \text{Prob\_content}_t + ((1 - \text{Prob\_content}_t) * \text{Learning}) \qquad (4)$$

where learning is kept constant at 0.01.

Then, the probability of judging beliefs on the basis of the source is decreased by the same amount according to equation 5:

$$\text{Prob\_source}_{t+1} = 1 - \text{Prob\_content}_{t+1} \qquad (5)$$

If the outcome does not satisfice, then the agent decreases the probability of using the strategy that produced this outcome and increases the probability of using the opposite strategy with future beliefs using equations 4 and 5. Note that both probabilities always add to 1.

If the individual decides not to try the belief, another belief is presented to the individual and the whole procedure is repeated. The agent keeps on trying beliefs for a specific number of time steps and the probability of using one or other strategy to judge beliefs may stabilize over time depending on the agents' competence and the environment's consistency.

## 2.2   Simulations and Data Collection

We ran simulations by varying the ability of the agent to judge the belief on the basis of its content or source. Hence, we assigned values to consistency and competence going from 0 to 1 by steps 0.01. By doing so, we ended up with a total of 10,100 combinations of values. All other parameters were kept constant. We ran 20 replications per combination of values and for each scenario representing the relationship between accuracy and utility. In total we ran 202,000 simulations per scenario. Simulations were stopped after the individual was presented with 1000, 3000 or 5000 beliefs, no matter whether the individual tried or not all the beliefs. For each combination of parameters and their 20 replications we calculated the percentage of replications where the agent tended to use the content strategy the most (i.e., probability of using content evaluation > 0.5) when evaluating the belief. The results presented here are the ones obtained after 1000 beliefs. Results for larger numbers of presented beliefs remain qualitatively the same as those of 1000 and are thus not presented.

## 3   Results

### 3.1   Scenario 1: Belief Utility Equals Belief Accuracy

Figure 1 presents the percentage of replications (n = 20) where content strategy was used by the agent, going from 0% (blue) to 100% (red) of the time. As expected, the agent used the content strategy the most when their competence was medium/high (i.e., > 0.5) and higher than the environment's consistency. On the contrary, when the environment's consistency was higher than the agent's competence and consistency was medium/high (i.e., > 0.5,) the agent relied more on the content strategy. When consistency and competence had a similar value (diagonal) or both of them had a medium/low value (i.e., < 0.5), the decision to use either of the strategies appears to be 50–50%.

### 3.2   Scenario 2: Belief Utility Equal to Uniform Distribution

The Fig. 2 presents the percentage of replications where content strategy was used by the agent, going from 0% (blue) to 100% (red) of the time. No clear pattern is observed, the prevalence of any of the strategies appears random for all combinations of consistency and competence values. Consistency and competence have a similar approximation value of 50% in most cases. We do not observe here any of the patterns seen in Figs. 1 and 3.

**Fig. 1** Heatmap of the percentage of replications [5] where the content strategy was used 100% (red) or 0% (blue) of the time for each combination of values of competence (x-axis) and consistency (y-axis)



**Fig. 2** Heatmap of the percentage of replications where the content strategy was used 100% (red) or 0% (blue) of the time for each combination of values of competence (x-axis) and consistency (y-axis)

### 3.3   Scenario 3: Belief Utility Equals 1—Belief Accuracy

Figure 3 again presents the percentage of replications where content strategy was used by the agent, going from 0% (blue) to 100% (red) of the time. In this case, we observe the reverse pattern of that from Fig. 1. The agent used the content strategy the most when its competence was lower than the environment's consistency and

**Fig. 3** Heatmap of the percentage of replications where the content strategy was used 100% (red) or 0% (blue) of the time for each combination of values of competence (x-axis) and consistency (y-axis)

when the environment's consistency was medium/high (i.e., > 0.5). On the contrary, the agent relied more on the source strategy when the agents' competence was higher than the environment's consistency and competence was medium/high (i.e., > 0.5). When consistency and competence had a similar value (diagonal) or both of them had a medium/low value (i.e., < 0.5), the decision to use either of the strategies appears to be 50–50%.

## 4 Discussion

In this model we have sought to understand how an agent can learn to change its epistemic vigilance strategy on the basis of two vital considerations. The first of these is the relation between its own competence in judging the content of the information that is presented and the degree to which the prima facie plausibility of the sources in its environment is consistently related to the accuracy of the information they present. The second is the relationship between the accuracy and utility of the presented information. What we have found is that the agent we modelled was able to flexibly modify its epistemic vigilance strategy depending upon both considerations. Thus, when presented with an environment in which utility and accuracy were closely connected and either consistency or competence were high, it was able to learn to primarily rely upon the epistemic vigilance strategy that could identify highly accurate beliefs, which (in that model) also had high utility. This scenario is an idealisation of most situations in which epistemic vigilance might be used, in that utility and accuracy typically are positively related. However, the inverse case was

particularly interesting to examine both because of the potential role that it could have in cases where beliefs function to motivate human cooperation, and because of the additional difficulty the agent faced within it. The second of these considerations is due to the fact that the inverse relationship between utility and accuracy meant that neither source nor content vigilance could reliably produce beliefs with high utility.

In this scenario, the best approach would be to avoid the epistemic vigilance strategy that was the more likely to lead to accurate beliefs and rely upon the other strategy to accept beliefs some of which, by accident, had high utility. Indeed, when presented with the scenario where utility and accuracy had an inverse relationship, we saw that the agent was able to avoid the epistemic vigilance strategy that was likely to generate highly accurate beliefs that, in that scenario, lacked in utility. Thus, the agent was able to react appropriately to the inverse relationship between accuracy and utility. Finally, we examined the scenario where there was no relationship between accuracy and utility of beliefs. In this case we also did not observe any relationship between the choice of epistemic strategy used and the levels of competence or consistency, showing that where there was no relationship between utility and accuracy, the agent had no particular tendency to prefer source over content vigilance.

This study showed that an agent can be capable of learning to alter its epistemic vigilance strategy in a way that reacts flexibly to its abilities, the epistemic structure of its environment and the relationship between the accuracy and utility of beliefs. However, it was quite limited in that it only considered these variables in very simple terms, did not explore other parameters, and did not look at the speed with which the agent was able to modify its behaviour to suit its conditions. In future work, we will explore these aspects of the studied phenomenon in order to achieve insight into how humans are capable of reacting flexibly to a range of situations.

In addition, in this model, only the cognitive processes of the agent were simulated, in order to determine regularities in its behaviour. The same variables have also been examined using a multi-agent model that allowed agent interaction to be focussed upon in the analysis [6].

# References

1. Sperber, D., Clément, F., Heintz, C., et al.: Epistemic vigilance. Mind Lang. **25**, 359–393 (2010). https://doi.org/10.1111/j.1468-0017.2010.01394.x
2. Talmont-Kaminski, K.: For god and country, not necessarily for truth: the nonalethic function of superempirical beliefs. Monist **96**, 447–461 (2013). https://doi.org/10.5840/monist201396320
3. Talmont-Kaminski, K.: Epistemic vigilance and the science/religion distinction. J. Cogn. Cult. **20**, 88–99 (2020). https://doi.org/10.1163/15685373-12340075
4. Simon, H.A.: The Sciences of the Artificial. Cambridge, M.I.T. Press (1969)

5.  Wickham, H.: ggplot2: Elegant Graphics for Data Analysis, 2nd edn. Springer, New York, NY (2016)
6.  Rybnik, M., Puga-Gonzalez, I., Shults, L., Dąbrowska-Prokopowska, E., Talmont-Kaminski, K.: An agent-based model of the role of epistemic vigilance in human cooperation (under review)

# A Simple Model of Citation Cartels: When Self-interest Strikes Science

**Davide Secchi** (ORCID)

**Abstract** This paper is an attempt to study a well known (probably little studied) phenomenon in academia: *citation cartels*. This is the tacit or explicit agreement among authors to cite each other more often than they would do in a more "sincere" approach to science. It can be intended as collusion and it can distort scientific progress in affecting a scholar's attention. The phenomenon has been around for decades and it does not seem to spare any discipline. By starting from outlining the characteristics of a "cartel," this study then builds an agent-based model in an attempt to define the extent to which colluding behavior affects progress in a given discipline by operating on citation counts. Data is still preliminary although enough to conclude that cartels promote lax scientific practices.

**Keywords** Citation cartels · Rigor · Scientific distortion · Agent-based modeling

## 1 Introduction

Citations are one of the currencies of modern scientists. They constitute the backbone of metrics that are used by hiring committees, grant-awarding institutions, promotion committees, and in general they are used to have a quick idea on someone's academic performance. Popular measures employed in the ranking of academic journals, such as the Impact Factor [10, 11], are calculated on citations that a journal attracts over a period of time. Other measures, such as Hirsch's index [15, 16], are also based on citations and are usually applied to individual scholars. More sophisticated measures, such as Scopus' CiteScore or the Eigenfactor Score, are still based on citations. In

---

---

D. Secchi (✉)
University of Southern Denmark, 4200 Slagelse, Denmark
e-mail: secchi@sdu.uk
URL: https://secchidavi.wixsite.com/dsweb

general, it seems difficult to dismiss the fact that citations play an important part in a scholar's career.

This chapter is concerned with those scholars that use citations as the goal of their scientific activities, mainly as a result of self-interest. In so doing, some enter a citation "agreement" (tacit or explicit) with other scholars in their community, aimed at citing each other's work more than that of others. This should give scholars who enter the agreement a citation advantage. Some have defined this phenomenon a *citation cartel* [5, 22]. The objective of this chapter is to explore the extent to which *citation cartels* distort the academic discourse and affect scientific progress. The analysis considers behavioral triggers of authors who engage in these 'cartels' and uses an agent-based simulation to explore its effects.

It is not the purpose of this chapter to discuss the appropriateness of using citations as a viable way to assess scientific work [3, 4, 13, 19]. This would entail a discussion of the purpose of citations, of whether, for example, they can be used as instruments of 'persuasion' [20]. Instead, the meaning of the initial part of this chapter is simply to remind readers that citations are used by academic institutions and peers to evaluate the quality of scholarly production. For this reason, the existence of *cartels*—i.e. tacit or explicit agreements that artificially steer scientific debates in a given direction— can be a problem for the progress of any discipline. Yet, the use of citations as the sole or as the main criterion to assess scholars and science is highly problematic. First, the idea that what is popular is inherently 'good' has much appeal, but it lacks a solid logical basis [30]. Second, citations can be ceremonial in the sense that an article can be included in a reference list just because it is customary to do so in certain fields (see, e.g., [8]). Third, the sheer number of citations does not tell anything about their use. Sometimes citations serve the scope of highlighting poor design of a study, theoretical and conceptual inconsistencies, or lack of rigor.

In spite of these concerns, citations are widely used. Hence, artificially boosting them casts a dark shadow over the way in which a discipline progresses. In the following, *citation cartels* are described and some examples are presented. An agent-based model is then introduced with a summary of preliminary findings. A short conclusion ends this chapter.

## 2   Citation Cartels

Although the phenomenon has been known for some time [9], there is not a wide literature on *citation cartels*. These have been defined as "groups of authors that cite each other disproportionately more than they do other groups of authors that work on the same subject" [7, abstract]. In the literature, "cartels" have been associated to journals [4, 14, 21] in which citation patters suggest a connection stronger than one may expect. In describing a borderline behavior, Davis [5] mentions the case of a review article published in *The Scientific World Journal* in which there were 124 references, 96 to the journal *Cell Transplantation* and 26 to the journal in which the review appeared. The review was published in 2010 and the citations to *Cell*

*Transplantation* were all from 2008 or 2009. The authors of the review article were editors of that journal. By repeating this behavior, the editors were able to boost the citation count for the journal and increase its Impact Factor.

This is an extreme case, specifically because the authors' behavior blatantly reveals the intentions behind it. In most cases, behavior is subtler and it is difficult to determine whether there is an open plan to boost citations or not. After all, science is arranged in communities of scholars that, among other actions, cite each other. This is ordinary behavior in academia, it serves the purpose of establishing a dialogue with colleagues with similar interests (the so-called "conversation"), it helps build consensus around a topic, and it makes it easier to understand key publications in the sub-field. For this reason, it is extremely difficult to determine whether a community of scholars engages in an actual "cartel" or not.

## 2.1 Behavioral Aspects of a Citation Cartel

There is more to extreme cases, though. "Cartel" is a word that implies there is collusion—the secret, typically unethical agreement among individuals aimed at artificially altering results. This may be too strong of a characterization for academic work.[1] But the nuances of the concept may be important to understand it better.

It is possible that a group of scholars agrees to cite each others more than they cite work outside of their scientific circle. This is something that can be done by either looking for other authors in the circle or by limiting the review of the literature to specific journals only—those where co-authors and others from that circle have agreed are the only one that count. The former behavior results in citations that are skewed towards a restricted number of authors, those who are also members of the 'circle.' The latter should also be visible in the list of references, where certain journals appear much more often than others.

The initial agreement may be a tacit understanding that scholars who share similar views (or a research agenda) would appear stronger if they cite each other's work. In this case, there is no explicit agreement, yet the effect on the sub-field is still that of a citation cartel. The effects of increasing citations are not limited to those in the cartel. In fact, it is a well-known fact in academia that high citations attract more citations. This is called "preferential attachment" and it is a distribution with power-law tails, based on the additional visibility that a highly cited article gains as opposed to other articles [1, 24]. Hence, the design of a cartel would make it such that highly cited papers (and this transfers to the authors) attract citations from scholars outside the circle. These are scholars who are diligent and address research questions more honestly, in an attempt to screen as many publications as possible on a given topic, without limiting themselves to any given list of journals or authors.

---

[1] Of course, this may depend on the field. Obviously, doctoring citation in areas of science where life and health of individuals are at stake may qualify more immediately as unethical.

For these scholars, the focus is the topic not the people (i.e. the circle). These have been called *sincere* scholars while the others have been called *strategic* scholars in a recent blog post [22].

## 2.2 Detecting a Cartel

The tools available to scientists to understand whether a set of citations qualify as a cartel are very limited. Fister and colleagues [7] presented an algorithm to detect cartels in a network. The essence of their logic is that authors who cite each other at a rate that is much higher than the average rate of citations in the network may have established a cartel. This is a good way to look at citation patterns however, it does not do enough to distinguish the existence of an actual (explicit or tacit) cartel from the emergence of a scientific interest. The point here is that those authors who cite each other often should be connected through means that are not easily attached solely to the science they produce. In other words, they may be part of the same academic society, work for the same organization, or have a past of being colleagues. Those are elements that may suggest there is an actual agreement.

Other measurements that can be used apply to journals rather than to authors [18]. This is not the focus of this chapter, but the logic involved is interesting. Kojaku et al. [18] use an algorithm for the detection of "anomalous citation groups" through the "CItation Donors and REcipients (CIDRE) algorithm." The idea is that of comparing communities to a null model—one that produces citations according to 'neutral' expectations of how science is supposed to evolve. This is particularly relevant although the problem is exactly that of isolating the 'neutral' model and its assumptions. Whatever is viable for one field may not be for another, as the citation patterns between disciplines clearly demonstrate.[2]

From the above, it is apparent that it is not easy to detect a cartel. Maybe, if the starting point is a discipline specific sub-field with a clear delimitation, patterns may be easier to study.

## 3 Modeling a Citation Cartel

The Citation Cartels simulation `CC1.1` is a fairly simple agent-based model. It has been developed with the aim of exploring to what extent scholars who engage in citation cartels are more likely to get ahead in the citation game and what consequences this "game" has on scientific work [6]. I could not find any other agent-based model

---

[2] A quick look at bibliometric indexes for top journals in different disciplines showcases the difference. For example, according to Journal Impact Factor by Clarivariate's Web of Science, in 2021 the journal *Lancet*—top journal for the field: 'general medicine'—was 202.731. The same index for the *Journal of Cultural Economy*—top journal for the field: 'cultural studies'—was 6.613.

of this phenomenon although author relations [12, 25, 26] and citations have been modeled before [17, 23].

The model has two agent types: (a) scholars, and (b) publications. They appear in an academic space that represents the environment in which these agents produce their scientific work. When they get closer in the model environment, they have more chances to cite each other's work. Proximity is, in this case, similar to the chance of knowing the work of others. Each step ($tick$) of the simulation is a year. The choice has been based on the fact that citations are usually counted by the various indexes with a yearly cadence. The simulation stops after 30 years.

### 3.1 Agent Characteristics

**Scholars**. Agent-scholars appear in the environment at random at the beginning of the simulation, and their number can be manipulated by the modeler $N_s[0, 1000]$. A proportion of them can be made to behave "politically" with the work they cite $P_s[0, 1]$. The remaining scholars do not have preferences in the way they cite, and follow a scientific logic. The former have been called *strategic* (*s*) and the latter *sincere* (*a*) scholars after Phelps [22].

Each agent-scholar screens the literature with some `rigor`, that is the extent to which one is willing to perform a wide search while writing a paper. The wider the search, the more sources a scholar may cite in their work. The modeler controls the top value of this characteristic ($max(r)$) and the values are assigned at random following a uniform distribution $U[0, 12]$.

Contrary to the *sincere*, the *strategic* agent-scholars can be disloyal to their circle—the community with which they engage in a cartel. No or minimal `deception` means that the only citations for these scholars are those coming from other *s*, while higher values of this variable indicate a more relaxed citation behavior. The modeler controls the mean `deception` ($\bar{d}$) and the values are attributed at random following a normal distribution with standard deviation $sd = 1$.

**Publications**. Agent-publications are generated by each scholar at random, with an exponential distribution around a floating point that can be chosen by the modeler ($P_e$). This choice follows the idea that very prolific authors are few while most publish at a rate that is far from those who peak. Publications when year $\geq 1$ are calculated as a random number selected between $[0, P_{c,0}]$, where $P_{c,0}$ is the publication count at the beginning to the simulation year $= 0$.

### 3.2 Procedures

In the system, connections represent citations. They are created by each agent-publication that searches the space around it, according to the formula:

$$\delta_{s_1,s_2} \leq r_{s_1} \times \left(1 + \frac{\max(d) - d_{s_1}}{\max(d)}\right) \qquad (1)$$

where $\delta_{s_1,s_2}$ is the distance between publication of strategic scholar $s_1$ and publication of strategic scholar $s_2$, $r_{s_1}$ is the rigor of $s_1$, $\max(d)$ is the maximum deception level in the system, and $d_{s_1}$ is the deception level for $s_1$. Higher numbers of $d_{s_1}$ increase the search range for publications of other strategic scholars. When $s$ is looking for $a$ to cite, then (1) changes sign:

$$\delta_{s_1,a_2} \leq r_{s_1} \times \left(1 - \frac{\max(d) - d_{s_1}}{\max(d)}\right) \qquad (2)$$

The search range is much closer to $s_1$ and it becomes more a function of $r_{s_1}$ rather than $d_{s_1}$. Yet, the lower deception values lead to a much reduced search range. Instead, sincere scholars $a$ only use their `rigor` to find other work publications to be cited.

In order to mimic the traction that highly cited work exercises on the academic community, each agent-scholar targets one highly cited work (at random) and moves towards it, such that it becomes increasingly likely that this work will be cited.

Another mechanism allows $s$ to turn $a$ and vice versa. When the mean citations to $a$'s work around $s$ is twice as much that of other $s$ around, then the $s$ turns. The same logic applies for $a$ but the mean citations for $s$ is only 20% superior. This is optional, it is called `adapt` and is controlled by the modeler.

## 4   Results

The simulation underwent a verification process [28] and a calibration process to test the parameter space [2]. Albeit simple, this simulation rapidly produces a number of links (citations) that spike up exponentially making calibration quite challenging. Due to this issue, a decision was made to present preliminary results. Further appropriate checks will be taken to integrate results in subsequent work.

I performed the simulation by selecting parameter values that were distant enough to guarantee an appropriate effect on the outcome variables. I kept the number of scholars constant at $N = 50$, as well as the input for the exponential distribution of publications $e = 1$. The other parameters assumed the following values: adapt $= \{T, F\}$, $P_s = \{0.2, 0.5, 0.8\}$, $d_s = \{1, 5, 9\}$, $\bar{r} = \{3, 6, 9\}$. Each run was repeated 180 times per each configuration of parameters (we calculated statistical power analysis as per [27, 29]).

Figure 1 shows the log number of citations per type of scholar at year $= 30$. In most cases when $P_s > 0.2$, the clouds of points are such that loyalty—that is here the opposite of deception and it happens when $\bar{d} = 1$—brings more references to $s$. As the mean value $\bar{d}$ increases and $s$ are 'free' to cite references outside of their circle, the numbers are such that the divide between the two types of scholars is not as wide. One way to interpret results in Fig. 1 is to use the diagonal $y = x$ as a reference

**Fig. 1** Number of citations for strategic and sincere scholars (year = 30, scholars = 50, $P_s$ = {0.2, 0.5, 0.8}, adapt = {$T$, $F$})

point. That is when $s$ and $a$ have the same number of citations; hence, points below the diagonal indicate more citations for $a$ while those above it show the opposite trend. An obvious trend is that $P_s$ attracts more citations when $s$ are the majority in the system. This is visible from the two panes at the right where $P_s = 0.8$. What is perhaps not that obvious is that this does not happen when they are a minority—panes on the left ($P_s = 0.2$). In this case the diagonal cuts across the points, with most being below it. The case in which the two types are split $P_s = 0.5$ clearly shows the effect of adapt, since the points are slightly towards $P_s$ when adapt = $F$ and always in favor of $P_s$ when adapt = $T$.

Figure 2 zooms in the case where $s$ start as a minority $P_s = 0.2$ but other scholars ($a$) have the possibility to switch if that is convenient enough. This figure shows LOESS curves with confidence intervals for each `deception` level, where $y$ is the mean citations for $s$ divided by the mean citations for $a$ ($\frac{c_s}{c_a}$) and $x$ is simply time. For low and mid-level of max $r$, the data seem to be not dispersed—i.e. the confidence intervals appear close to the line. A low level of deception ($d = 1$) seem to give $s$ quite a substantial advantage over $a$, with $s$ citations that are, on average, two or three times higher than those of $a$. This clearly explains why many $a$ switch to become $s$, following self-interest. The remaining two curves in the left and center pane overlap. There is still a consistent advantage for $s$ (this is the condition of this selection) but $d$ does not seem to explain it, while $\bar{r}$ provides a better explanation.

The pane on the right of Fig. 2 is different. In fact, it seems that max $r$ explains the $s$ advantage while $d$ does not. In this case, a simple rigorous academic work is enough for these scholars to get an advantage, independent of whether they are loyal to their circle or not. The increasing number of $s$ scholars guarantees them that behaving strategically pays off only because they are a majority in the system. The counterargument is offered by Fig. 3. The logic is the same as Fig. 2 but initial proportion $P_s = 0.2$ is not allowed to change, i.e. adapt = $F$. This is the case where

**Fig. 2** LOESS regression curves and confidence intervals for $\bar{c}_s/\bar{c}_a$ over time (scholars = 50, $P_s = 0.2$, adapt = $T$, $N_s > 25$)



**Fig. 3** LOESS regression curves and confidence intervals for $\bar{c}_s/\bar{c}_a$ over time (scholars = 50, $P_s = 0.2$, adapt = $F$)

scholars stick to their beliefs on how to conduct scientific research. The effect of $d$ is similar to that observed in Fig. 2. However, in this case the curve is rather different. In fact, it grows for the first half of the simulation time ($\approx 12$ or 15 years) and then declines. Rigor, in this case, reduces the advantage of strategic behavior, especially when the overall number of publications increase and $s$ remain stuck on limited numbers. Deception seems to help, but it is irrelevant when rigor increases (pane on the right).

## 5 Concluding Remarks

The simulation demonstrates that strategic behavior and cartels do give a citation advantage to those who practice it. Also, when we assume that all academics are selfish- i.e. 'play the game' - and want to gain from this behavior, this advantage ceases to exist since the vast majority of academics switch to a strategic behavior. Indeed, this is a strong assumption to make over a scholar's behavior. Instead, in a system where there are stable preferences over scientific practices, the advantage of $s$ scholars is limited in time and seems to fade away as a field progresses such that more rigorous practices become more common. In this case, the consequence is that the advantage of cartels remains temporary. However the assumption of stable preferences is probably too strict.

Another learning point is that, to exploit the cartel in full, strategic scholars need to be rather "purists" and loyal to their circle. This would, on average and when rigorous search is not a priority, increase their advantage. In other words, lax scientific standards (low rigor) are rewarded with a citation increase, rather than punished as bad science.

Overall, not only are citation cartels unethical and capable of steering a system in the direction of the publications coming from its members, but they promote less rigorous scientific practices that, in the long run, may affect results and relevance of a scientific field. This is how the self-interested scholar damages the scientific enterprise.

## References

1. Barabási, A.L., Albert, R.: Emergence of scaling in random networks. Science **286**(5439), 509–512 (1999)
2. Boero, R., Squazzoni, F.: Does empirical embeddedness matter? Methodological issues on agent-based models for analytical social science. J. Artif. Soc. Soc. Simul. **8**(4), 6 (2005)
3. Callaway, E.: Beat it, Impact Factor! Publishing elite turns against controversial metric. Nature **535**(7611), 210–211 (2016)
4. Chakraborty, J., Pradhan, D., Nandi, S.: Research misconduct and citation gaming: a critical review on characterization and recent trends of research manipulation. Lecture Notes on Data Engineering and Communications Technologies, vol. 71, pp. 485–492 (2022). https://doi.org/10.1007/978-981-16-2937-2_30, cited by 1
5. Davis, P.: The emergence of a citation cartel. The Scholarly Kitchen (blog), vol. 15 (2012)
6. Edmonds, B., Page, C.L., Bithell, M., Chattoe-Brown, E., Grimm, V., Meyer, R., Montanola-Sales, C., Ormerod, P., Root, H., Squazzoni, F.: Different modelling purposes. J. Artif. Soc. Soc. Simul. **22**(3), 6 (2019)
7. Fister Jr., I., Fister, I., Perc, M.: Toward the discovery of citation cartels in citation networks. Front. Phys. 49 (2016)
8. Foss, N.J.: Bounded rationality in the economics of organization: "much cited and little used". J. Econ. Psychol. **24**(2), 245–264 (2003)
9. Franck, G.: Scientific communication—a vanity fair? Science **286**(5437), 53–55 (1999)
10. Garfield, E.: Journal impact factor: a brief review. Can. Med. Assoc. J. **161**(8), 979–980 (1999)
11. Garfield, E.: The meaning of the Impact Factor. Int. J. Clin. Health Psychol. **3**(2), 363–369 (2003)

12. Gilbert, N.: A simulation of the structure of academic science. Sociol. Res. (online) **2**(2), 91–105 (1997)
13. Haley, M.: On the inauspicious incentives of the scholar-level h-index: an economist's take on collusive and coercive citation. Appl. Econ. Lett. **24**(2), 85–89 (2017). https://doi.org/10.1080/13504851.2016.1164812, cited by 24
14. Hickman, C., Fong, E., Wilhite, A., Lee, Y.: Academic misconduct and criminal liability: manipulating academic journal impact factors. Sci. Public Policy **46**(5), 661–667 (2019). https://doi.org/10.1093/scipol/scz019, cited by 8
15. Hirsch, J.E.: An index to quantify an individual's scientific research output. Proc. Natl. Acad. Sci. **102**(46), 16569–16572 (2005)
16. Hirsch, J.E.: Does the h index have predictive power? Proc. Natl. Acad. Sci. **104**(49), 19193–19198 (2007)
17. Ionescu, G., Chopard, B.: An agent-based model for the bibliometric h-index. Eur. Phys. J. B **86**, 1–7 (2013)
18. Kojaku, S., Livan, G., Masuda, N.: Detecting anomalous citation groups in journal networks. Sci. Rep. **11**(1) (2021). https://doi.org/10.1038/s41598-021-93572-3, cited by 3
19. Moed, H.F., Van Leeuwen, T.N.: Impact factors can mislead. Nature **381**(6579) (1996)
20. Gilbert, N.: Referencing as persuasion. Soc. Stud. Sci. **7**(1), 113–122 (1977)
21. Perez, O., Bar-Ilan, J., Cohen, R., Schreiber, N.: The network of law reviews: citation cartels, scientific communities, and journal rankings. Mod. Law Rev. **82**(2), 240–268 (2019). https://doi.org/10.1111/1468-2230.12405, cited by 5
22. Phelps, R.P.: How citation cartels give 'strategic scholars' an advantage: a simple model. Retraction Watch (blog) (2022)
23. Pluchino, A., Burgio, G., Rapisarda, A., Biondo, A.E., Pulvirenti, A., Ferro, A., Giorgino, T.: Exploring the role of interdisciplinarity in physics: success, talent and luck. PLoS ONE **14**(6), e0218793 (2019)
24. Price, D.J.d.S.: A general theory of bibliometric other cumulative advantage distributions. J. Am. Soc. Inf. Sci. **27**, 292–306 (1976)
25. Secchi, D., Cowley, S.J.: Modeling organizational cognition: the case of impact factor. J. Artif. Soc. Soc. Simul. **21**(1), 13 (2018)
26. Secchi, D., Cowley, S.J.: Improbable fairness: reviewing under the lenses of Impact Factor. RASK **50**(Autumn), 191–209 (2019)
27. Secchi, D., Seri, R.: Controlling for 'false negatives' in agent-based models: a review of power analysis in organizational research. Comput. Math. Organ. Theory **23**(1), 94–121 (2017)
28. Seri, R., Martinoli, M., Secchi, D., Centorrino, S.: Model calibration and validation via confidence sets. Econometr. Stat. **20**(October), 62–86 (2021)
29. Seri, R., Secchi, D.: How many times should one run a computational simulation? In: Edmonds, B., Meyer, R. (eds.) Simulating Social Complexity. A Handbook, 2nd edn., pp. 229–251. Springer, Heidelberg (2017)
30. Woods, J.: The Death of Argument: Fallacies in Agent-Based Reasoning. Kluwer, Dordrecht (2004)

# A Study on Multi-scale Modeling in Social Simulation Focusing on Relationships Among Decision-Makers

**Gaku Shimizu, Toshiya Kaihara, Daisuke Kokuryo, and Nobutada Fujii**

**Abstract**  The concept of System of Systems (SoS) is important to realize a society that creates sustainable value through the coordination and cooperation of systems. However, for SoS consisting of subsystems at different spatio-temporal levels, conventional modeling methods are applied independently at each level, making it impossible to conduct both macro and micro evaluations at the same time. In this paper, we propose a multi-scale modeling method that enables modeling of each system component at different spatio-temporal levels. The proposed method is applied to a local city under COVID-19, and a comprehensive analysis of the target system is conducted by modeling and integrating decision-makers at different levels: citizens, organizations, and municipality.

## 1  Introduction

Society 5.0 [1] is a concept proposed by Japanese Government in 2016, which creates sustainable value and services through the coordination and cooperation of systems. The concept of System of Systems (SoS) is important to realize Society 5.0 [2]. SoS is a collection of elements, each of which can be considered a system, with operational independence in that each elemental system can operate independently,

G. Shimizu (✉) · T. Kaihara · D. Kokuryo · N. Fujii
Graduate School of System Informatics, Kobe University, 1-1 Rokkodai-Cho, Nada, Kobe, Hyogo 657-8501, Japan
e-mail: shimizug@kaede.cs.kobe-u.ac.jp

T. Kaihara
e-mail: kaihara@kobe-u.ac.jp

D. Kokuryo
e-mail: kokuryo@port.kobe-u.ac.jp

N. Fujii
e-mail: nfujii@phoenix.kobe-u.ac.jp

33

**Fig. 1** Concept of system of systems



and managerial independence in that each elemental system has the authority to manage the system [3]. One of the characteristics of SoS is that it includes various subsystems at different spatio-temporal levels. These systems have a hierarchical structure, mediated by the macro-level of the society as a whole, and are mixed with various heterogeneous systems with lower levels at the meso- and micro-levels (Fig. 1).

On the other hand, it is difficult to achieve both a macro-level evaluation of the entire system and a micro-level evaluation of the details simultaneously, because conventional modeling methods are applied independently at each level, and there is a need for a modeling method that can appropriately design and analyze social systems with SoS structures. Recently, there has been much research on hybrid simulation (HS: defined as a modeling approach that combines two or more of the following methods: Discrete-Event Model (DEM), System Dynamics (SD), and Agent-Based Model (ABM)) to represent complex systems [4], but few studies have explicitly considered differences in spatio-temporal level of the components. Therefore, we focus on multi-scale modeling, in which models are made by using modeling methods appropriate for each level and integrated, and propose a modeling method that can efficiently implement models with the accuracy required at each level and enable evaluation suited to each level. The proposed method is applied to a local city where COVID-19 is prevalent, and the efficiency is evaluated with simulation experiments.

## 2 Multi-scale Modeling in Social Simulation

In this chapter, multi-scale modeling in social simulation is proposed. First, we describe the concept of multi-scale modeling. Then, we describe the systematization of decision-makers and the methods for integrating multiple decision-making models, which are necessary for modeling social systems by multi-scale modeling.

**Fig. 2** Comparison of conventional method and multi-scale modeling

## 2.1 Multi-scale Modeling

Real-world phenomena span a wide range of temporal and spatial scales. Conventional modeling methods that focus on one scale cannot capture phenomena near the boundaries of these scales because they require assumptions and conditions for other scales (Fig. 2: Left). It is also difficult to achieve both a macro-level evaluation of the entire system and a micro-level evaluation of the details simultaneously from the viewpoints of efficiency and accuracy. Therefore, we have focused on multi-scale modeling, which can simultaneously consider models at different levels (Fig. 2: Right). This concept is considered to be useful in several fields, because it can simultaneously capture changes at different spatio-temporal levels by seamlessly connecting diverse systems that include different scales [5].

For social systems consisting of various decision-makers and having a hierarchy, it is also necessary to seamlessly connect micro-models such as individual behavior models, meso-models such as industrial structure models, and macro-models such as policy and economic evaluation models, and simultaneously analyze them. In the social systems, the decision-makers at different spatio-temporal levels are regarded as elemental systems of SoS, and then they can be seamlessly integrated. By realizing multi-scale modeling, it is possible to implement models at the different spatio-temporal levels with accuracy and efficiency for the target social system. Then, the model enables analysis suitable for the entire target system and each level. On the other hand, to apply multi-scale modeling into social systems, it is necessary to systematize decision-makers at different spatio-temporal levels and to integrate multiple decision-making models.

## 2.2 Systematization of decision-makers

To apply multi-scale modeling into social systems, the framework proposed by DeLaurentis et al. [6] is applied to organize the system components in this paper. The systematization procedure is as follows.

**Step 1**: Classify system components into the following four types.

- Resources: Physical components of the system
- Operations: Activity policies and procedures for physical components
- Economics: Nonphysical components in a market economy
- Policy: Policies affecting each component

**Step 2**: Classify components within each category into several hierarchies based on their relative position at the spatio-temporal level.

The existence of decision-makers at various levels can be grasped by focusing on Resources elements.

## 2.3 Modeling of Decision-Makers and Their Integration

We construct an ABM consisting of entities that are systematized by the method described in Sect. 2.2 and have a decision-making model appropriate for that level. When the decision-making models are at different spatio-temporal levels, the intervals between decisions and the range of their influence are different. Therefore, it is necessary to integrate the models seamlessly, considering the spatio-temporal level and the relationships among decision-makers, rather than simply connecting inputs and outputs among the models.

In this paper, the target society at the local city scale is divided as follows according to the framework in Sect. 2.2: micro-model for citizens at $\alpha$-level, meso-model for organizations at $\beta$-level, and macro-model for a municipality at $\gamma$-level. Then, by focusing on the policy structure of the municipality, we find two types of relationships among these entities, and propose a method to seamlessly integrate the three models by explicitly introducing the concept of "social consciousness" into the model.

**Constructing Relationships among Decision-Makers based on Policy Structure**. In this paper, we construct relationships among citizens, organizations, and a municipality based on policy structure. Policies have a hierarchical structure consisting of a chain of objectives and means, and are generally classified into the categories of "Policy" (in a narrow sense), "Program", and "Project" [7]. The higher the hierarchy, the more abstract, and the lower the hierarchy, the more concrete. The characteristics of each hierarchy are as follows.

- **Policy**: Major set of administrative activities aimed at a major goal.
- **Program**: Set of administrative activities aimed at a concrete objective based on the major goal.
- **Project**: Basic unit of administrative activity to implement programs.

Based on the above characteristics, two types of relationships between decision-makers are established. The role of policy is to encourage citizens and organizations to change their behavior by appealing to their social consciousness. Therefore, the three decision-makers have a relationship through "Social Consciousness". On the

other hand, the role of project is to change the environment for citizens' activities by imposing rules and institutions on organizations. Therefore, the three decision-makers have a relationship through "Code of Conduct".

**Introduction of Social Consciousness**. It is inspired by homogenization used in multi-scale modeling in materials science. Homogenization is a method of replacing non-homogeneous material with periodic microstructures with macroscopically equivalent homogeneous material [8]. In a social system, this is regarded as the process of extracting the consciousness shared by individual entities in society. In this paper, social consciousness is defined as "the homogenized consciousness of the entire society generated by the decisions and actions of the decision-makers in the target system". The seamless integration of the three decision-making models is achieved by modeling the social consciousness shared by decision-makers at different levels.

**Integration Based on Two Types of Relationships**. Based on the two types of relationships and modeling of social consciousness, we seamlessly integrate the decision-making models at different levels. The model overview is shown in Fig. 3, and the details of the model integration are described below.

*Social Consciousness (Fig. 3: Right)*. The output of micro/meso-model, i.e., the results of the actions of citizens and organizations, generates social consciousness, which becomes the input of macro-model and influences the policy decisions of the municipality. On the other hand, social consciousness is controlled by the policies of the municipality, which is the output of macro-model and influences the decision-making of citizens and organizations in micro/meso-model.

*Code of Conduct (Fig. 3: Left).* The output of macro-model, i.e., the projects of the municipality, imposes rules and institutions on the organizations as inputs to meso-model. In response, organizations follow or go against the rules and institutions according to their management policies. In meso-model, changes in the decision-making of the organization cause changes in the activity environment of citizens, and citizens in micro-model act with the changes.

## 3 Application to the COVID-19 Problem

The proposed method described in the previous section is applied to a social system assuming a local Japanese city where COVID-19 is prevalent. In this city, there are various decision-makers at different spatio-temporal levels who make decisions on measures to control and prevent the spread of the disease. We model the target city in detail using synthetic population data [9] that estimates the national municipal population using the household composition restoration method. However, since we focus on decision-makers within the target city, we assume that there is no population outflow or inflow from the target city.

**Fig. 3** Proposed model with social consciousness

**Table 1** Classification of organizing the system components

| Level | Resources | Operations | Economics | Policy |
|---|---|---|---|---|
| $\gamma$ | Municipality | Infection control measures, vaccination projects | Budget, project expenses | – |
| $\beta$ | Office, store, school, hospital, vaccination venue | Organization operation, inoculation services | Budget, profit, operating expenses | Reduced business hours, school closure, securing hospital beds, vaccine distribution |
| $\alpha$ | Citizen, bed, vaccine, medicine | Citizen behavior, management of medical supplies | Labor cost, medical expenses | Restrictions on going out, telework, vaccination |

## 3.1 Systematization of Decision-Makers

Table 1 shows the classifications of organizing the system components using the framework described in Sect. 2.2 for the target system. In this table, the spatio-temporal influence of the system components increases in the order of $\alpha$, $\beta$, and $\gamma$. By focusing on Resources in Table 1, the existence of decision-makers at various levels, i.e., citizens at $\alpha$-level, organizations at $\beta$-level, and a municipality at $\gamma$-level, is grasped.

## 3.2 Modeling of Decision-Makers and Their Integration

Using the method described in Sect. 2.3, each systematized decision-maker is modeled and integrated as multi-scale modeling. Citizens at $\alpha$-level are micro-entities, characterized by generating emergent phenomena bottom-up through their interac-

tions. Therefore, they are modeled as micro-model. The decision-makers of offices, stores, schools, and vaccination venues among $\beta$-level organizations are modeled as meso-model, which are the activity environments of citizens. The municipality at $\gamma$-level is characterized by its decision-making based on macro variables and its top-down influence on the social system. Therefore, it is modeled as macro-model.

**Micro-model**. We model citizens' decision-making and their interactions as micro-model. It is assumed that citizens live in the local Japanese city consisting of multiple areas where households and $\beta$-level organizations exist.

*Citizen Type*. Citizens are classified into the following four types by employment status and age among the attributes of the synthetic population data.

– **Regular employees**: Employment status is general worker.
– **Non-regular employees**: Employment status is short-time worker or temporary worker and over 19 years old.
– **Students**: Employment status is not general worker, 7–18 years old.
– **Others**: Other than above.

*Citizen Action*. The daily action flow of citizens is shown in Fig. 4. The citizens determine destinations if they have no infection symptoms. The destinations that citizens can go to are defined by their attributes. The product of the outing rate *or* and the place-specific rate *r* determines whether or not citizens go to each place. Citizens with more than one destination go out and spend *c* yen per visit to a store. Citizens visit and stay at all destinations, and then return home.

*Disease Transition and Transmission*. Based on the literature on agent simulation for COVID-19 [10], a citizen has the seven states shown in Fig. 5. The citizen has infectivity from two days before transitioning to Mild or Asymptomatic to the Recovered, and spreads the infection with a defined probability (= contact rate $cr *$ transmission rate $tr$) at home and the destination. As shown in Fig. 4, agents with symptoms (circled by the dotted line in Fig. 5) stay in their house, and are home-cured or quarantined (no infection transmission) according to probability $hr$.

*Vaccination*. This paper assumes that the vaccine is highly effective, and 90% of infections can be controlled by the vaccine. Then, citizens can reduce the incidence $p_{sym}$ and severity $p_{sev}$ of disease through vaccination.

**Meso-model**. Organizations are modeled as activity environments for citizens. In this paper, since the target problem is COVID-19 emergency, organizations follow the instructions of the municipality and don't make any decisions, such as management policy decisions. Each organization has the following features.

– **Office**: Workplace for both regular and non-regular employees.
– **Store**: A place where citizens do consumption activities.
– **School**: Elementary, junior high, and high schools where students are assigned according to their address and age.
– **Vaccination venue**: A place where vaccines are allocated by the municipality and where citizens receive vaccines.

**Fig. 4** Action flow of citizen



**Fig. 5** Disease transition



**Macro-model**. The decision-making of a municipality is modeled as macro-model. Based on the policy structure described in Sect. 2.3, the municipality has policies and projects. Each of them is described as follows.

*Policy*. We introduce the concept of stages in the model, based on the Japanese government's COVID-19 measures. These stages represent the level of reinforcement of COVID-19 measures by the municipality and are classified into four stages according to defined criteria. The following situations are assumed in each stage, and a set of projects to be implemented in each stage is prepared.

- **STAGE1**: No measures are required.
- **STAGE2**: Caution should be strengthened.
- **STAGE3**: Measures should be strengthened.
- **STAGE4**: Maximum measures are required.

**Project**. Projects are a means of achieving the policies and we introduce them into the model as specific measures for organizations. In this paper, the municipality conducts a project to request reduced business hours and provide vaccines.

**Integration of the Three Models**. We integrate the three models according to the model overview shown in Fig. 3. In this paper, however, organizations are assumed to fully accept projects from the municipality in response to the COVID-19 emergency, and no connection to social consciousness is considered. The model integration by decision-maker relationship is described below.

*Model Integration by Social Consciousness*. In this paper, social consciousness is assumed to be formed based on the number of infected people output from micro-model as a result of the actions of citizens and modeled by SD (Fig. 6). Table 2 shows the definitions of the symbols and the formulation is as follows.



**Fig. 6** Stock and flow diagram

**Table 2** Definition of symbols

| $C$ | Social consciousness | $dr$ | Decrease rate of social consciousness |
|---|---|---|---|
| $IN$ | Inflow of $C$ | $NI$ | Number of infected people (from micro-model) |
| $OUT$ | Outflow of $C$ | $ni_0$ | Standard for the number of infected people |
| $CI$ | Indicators of social consciousness | $PV_{in}$ | Policy variable for inflow (from macro-model) |
| $lag$ | Lag in the generation of social consciousness | $PV_{out}$ | Policy variable for outflow (from macro-model) |
| $k$ | Time discount rate | $t$ | Simulation time (unit: day) |

$$\frac{dC(t)}{dt} = IN(t - lag) - OUT(t) \tag{1}$$

$$OUT(t) = C(t) * dr * (1 + PV_{out}(t)) \tag{2}$$

$$C(0) = 0 \tag{3}$$

$$IN(t) = \frac{CI(t)}{1 + k * t} * (1 + PV_{in}(t)) \tag{4}$$

$$CI(t) = \frac{NI(t)}{ni_0} \tag{5}$$

Equation (1) represents the change in $C$ per unit time, and $IN$ is the value before $lag$ day. Equation (4) represents that $IN$ is determined by $CI$ and decreases with time. Equation (2) shows that $OUT$ is a constant multiple of $C$. Equation (5) shows that the indicator of social consciousness is determined by the number of infected people. Equation (3) shows that the initial value of $C$ is 0.

The model integration is performed via the social consciousness formulated above. The municipality in macro-model makes policy decisions using $C$ as an input; the stages change according to the value of $C$, and $PV_{in}$ and $PV_{out}$ determined by each stage control the value of flow as represented by (4) and (2). On the other hand, citizens in micro-model make decisions by $C$, which changes according to (1), and $cs$, a parameter that represents the citizen's individual sense of crisis about infectious diseases. Specifically, in the "Citizen Action" item, the outing rate *or* of citizens is multiplied by $1 - C * cs$. In the "Vaccination" item, if $C * cs$ exceeds the threshold value $T$ for vaccination, which is given to all citizens in common, the citizen tries to vaccinate.

**Model Integration by Code of Conduct**. The project according to the stage determined by macro-model is implemented for the organization that is meso-model. This changes the activity environment of the citizens in micro-model. Specifically, in the parts described in meso-model and micro-model items, the time when citizens can go to the store becomes shorter because the closing time of the store is earlier due to the request for reduced business hours. In addition, citizens who satisfy all of the following conditions in the vaccination project: target age, vaccine inventory, and $C * cs$ exceeding the $T$, make a reservation for vaccination, and the reserved amount is allocated to the vaccination venue.

## 4  Computational Experiments

To analyze the behavior of the target system described in the previous chapter, the social simulation is performed using AnyLogic [11], a modeling tool that enables simulation by combining three modeling methods (SD, DEM, and ABM) within a single model. The target city is Ashiya City, a small local city in Japan.

## 4.1 Experimental Conditions

The experimental conditions are shown in Table 3 (basic condition), Table 4 (parameters by age), and Table 5 (probability of visit by type). The conditions related to the target city ($Q$, $R$, $P$, $W$, $S$, $E$, $J$, $H$, $V$) are set from the statistical [12] and synthetic population data [9] of Ashiya City. The conditions related to COVID-19 symptoms ($ET$, $MST$, $MRT$, $ART$, $SRT$, $SDT$, $p_{sym}$, $p_{sev}$, $p_{dea}$) are set based on the data from the Japanese Ministry of Health, Labor and Welfare [13] and literature on agent simulation for COVID-19 [10]. The conditions for citizen actions ($r_{of}$, $r_{st}$, $r_{sc}$) are set based on data from the Japanese Ministry of Land, Infrastructure, Transport and Tourism [14]. The sense of crisis $cs$ is set based on the assumption that it becomes stronger with increasing age. The other conditions are set through a preliminary experiment.

Table 6 shows the parameters for policies and projects that change according to $C$. The municipality determines the stage transition once every 4 weeks, and the stage is assumed to transition one stage at a time, even if the value of $C$ changes rapidly. The vaccination schedule is as follows: 26 weeks after the date of reaching STAGE2, vaccination becomes available for those aged 65 years and older, followed by those aged 18 years and older, and then those aged 12 years and older, in order of 8 weeks intervals. The number of simulation trials is 100.

**Table 3** Basic experimental conditions

| | | | | | |
|---|---|---|---|---|---|
| Area number $Q$ | 8 | Consumption amount $c$ | [1000, 10,000] | Exposed duration $ET$ | $N(5, 1)$ |
| Household number $R$ | 40,102 | Home care rate $hr$ | 0.7 | Mild-severe duration $MST$ | $N(7, 1)$ |
| Population $P$ | 90,244 | Vaccine threshold $T$ | 0.4 | Mild-Recovered duration $MRT$ | $N(10, 1)$ |
| Office number $W$ | 3111 | Transmission rate with symptoms $tr_{sym}$ | 0.1 | Asymptomatic-recovered duration $ART$ | $N(10, 1)$ |
| Store number $S$ | 1530 | Transmission rate without symptoms $tr_{asy}$ | 0.05 | Severe-recovered duration $SRT$ | $N(7, 1)$ |
| Elementary school number $E$ | 8 | Household contact rate $cr_{hh}$ | 0.4 | Severe-dead duration $SDT$ | $N(7, 1)$ |
| Junior high school number $J$ | 3 | Office contact rate $cr_{of}$ | 0.07 | Lag $lag$ | 5 |
| High school number $H$ | 3 | Store contact rate $cr_{st}$ | 0.09 | Time discount rate $k$ | 0.005 |
| Vaccination venue number $V$ | 1 | School contact rate $cr_{sc}$ | 0.01 | Decrease rate $dr$ | 0.05 |
| Outing rate standard $or_0$ | 1 | Initial exposed number $IE$ | 20 | Infected people standard $ni_0$ | 25 |

**Table 4** Parameters by age

| Age | Incidence, $p_{sym}$ | Severity, $p_{sev}$ | Lethality, $p_{dea}$ | Crisis sense, $cs$ |
|---|---|---|---|---|
| 0–9 | 0.50 | 0.00050 | 0.00003 | $N(12.0, 1)^{-1}$ |
| 10–19 | 0.55 | 0.00165 | 0.00008 | $N(11.5, 1)^{-1}$ |
| 20–29 | 0.60 | 0.00720 | 0.00036 | $N(11.0, 1)^{-1}$ |
| 30–39 | 0.65 | 0.02080 | 0.00104 | $N(10.5, 1)^{-1}$ |
| 40–49 | 0.70 | 0.03430 | 0.00216 | $N(10.0, 1)^{-1}$ |
| 50–59 | 0.75 | 0.07650 | 0.00933 | $N(9.5, 1)^{-1}$ |
| 60–69 | 0.80 | 0.13280 | 0.03639 | $N(9.0, 1)^{-1}$ |
| 70–79 | 0.85 | 0.20655 | 0.08923 | $N(8.5, 1)^{-1}$ |
| 80+ | 0.90 | 0.24570 | 0.17420 | $N(8.0, 1)^{-1}$ |

**Table 5** Probability of visit by type

| Type | Office, $r_{of}$ | Store, $r_{st}$ | School, $r_{sc}$ |
|---|---|---|---|
| Regular | 0.77 | 0.18 | 0.00 |
| Non-regular | 0.62 | 0.29 | 0.00 |
| Students | 0.00 | 0.18 | 1.00 |
| Others | 0.00 | 0.52 | 0.00 |

**Table 6** Relationship between social consciousness and policy/project parameters

| Social consciousness, $C$ | STAGE | $PV_{in}$ | $PV_{out}$ | Closing time, $ct$ | Vaccine supply, $v$ |
|---|---|---|---|---|---|
| 0.0–1.5 | 1 | 0.2 | 0.8 | 24 | 1000 |
| 1.5–3.0 | 2 | 0.4 | 0.6 | 24 | 2000 |
| 3.0–4.5 | 3 | 0.6 | 0.4 | 20 | 3000 |
| 4.5+ | 4 | 0.8 | 0.2 | 20 | 4000 |

## 4.2 Experimental Results

We simulated 730 days in the local Japanese city under the experimental conditions described in the previous section. The mean of the total number of infected people was 8280.7 ± 1632.1 in 100 trials. The real total number of infected people in Ashiya City for two years (April 2020 to March 2022) was 5403, and the simulation results are valid for the scale of infection, taking into account the omission of tests. Since the results of each trial showed similar trends, the results of one trial will be discussed below.

Figure 7 shows the number of people in each state, Fig. 8 shows the social consciousness and stage, Fig. 9 shows the consumption amount, and Fig. 10 shows the number of completed vaccinations. First, we note that in Fig. 7, the wave of the num-

Fig. 7 Number of people in each state

Fig. 8 Social consciousness and STAGE

Fig. 9 Consumption amount

ber of infected people increases and decreases repeatedly. As shown in Fig. 8, as the number of infected people increases, *C* increases, and the outing rate of citizens decreases, resulting in a decrease in the number of infected people. The reverse is also true, and thus the number of infected people vibrates up and down. In addition, as shown in (4), social consciousness is less likely to increase with time, and the wave of the number of infected people becomes larger after 400 days.

**Fig. 10** Number of completed vaccinations



Next, we focus on the relationship between social consciousness and stage in Fig. 8. Figure 8 shows that municipalities raise or lower their stage by increases or decreases in $C$. Social consciousness changes every day, while the municipality makes a decision once every four weeks. Then the municipality is not able to take prompt action, and the timing of the peak of $C$ and STAGE4 is not aligned.

Finally, in Figs. 9 and 10, we focus on the relationship between each project of the municipality and social consciousness and stage. Figure 9 shows that (a) the consumption amount increases or decreases continuously and (b) it increases or decreases discontinuously. (a) is because the outing rate increases or decreases with the change in social consciousness. (b) occurs at the border between STAGE2 and 3 in Fig. 8. This is because the closing time of stores is uniformly moved from 24:00 to 20:00 after STAGE3 due to the municipality's request for reduced business hours. Thus, it can be said that citizens' actions change in response to changes in both their social consciousness and the environment of their activities. Moreover, Fig. 10 shows that citizens vaccinate the elderly first in response to the municipality's vaccination project. In particular, the number of vaccinations increases as the threshold for vaccination is exceeded at the time when social consciousness increases. Figure 7 shows that as the percentage of completed vaccinations increases, the number of people with symptoms is reduced.

Thus, the proposed method enables analysis focusing on decision-makers at each level. Furthermore, compared to previous studies on COVID-19 (e.g., [10]), the proposed method enables not only analysis of results for specific measures or situations, but also a comprehensive analysis of the entire social system.

## 5   Conclusion

In this paper, we proposed a multi-scale modeling with social consciousness in social simulation. The proposed method was applied to a local city under COVID-19, and the results from the social simulation showed that the proposed method is capable

of performing a comprehensive analysis of the target system. In future works, the decision-making algorithm for each decision-maker should be sophisticated, and proposals for institutional design of social systems that reflect the intentions of different levels of decision-makers should be made.

# References

1. Society 5.0, https://www8.cao.go.jp/cstp/english/society5_0/index.html. Accessed 08 Feb 2023
2. Kaihara, T., et al.: Innovative Systems Approach for Designing Smarter World. Springer (2021)
3. Maier, M.W.: Architecting principles for systems-of-systems. Syst. Eng. **1**(4), 267–284 (1999)
4. Brailsford, S.C., et al.: Hybrid simulation modelling in operational research: a state-of-the-art review. Eur. J. Oper. Res. **278**(3), 721–737 (2019)
5. Weinan, E.: Principles of Multiscale Modeling. Cambridge University Press (2011)
6. DeLaurentis, D., et al.: A system-of-systems perspective for public policy decision. Rev. Policy Res. **21**(6), 829–837 (2004)
7. Policy Evaluation Implementation Guidelines, https://www.soumu.go.jp/main_content/000556222.pdf. Accessed 08 Feb 2023
8. Terada, K., et al.: A class of general algorithms for multi-scale analyses of heterogeneous media. Comput. Methods Appl. Mech. Eng. **190**(40–41), 5427–5464 (2001)
9. Harada, T., et al.: Projecting synthetic households on buildings using fundamental geospatial data. In: Proceedings of Social Simulation Conference 2017, 10 pp. (2017)
10. Kerr, C.C., et al.: Covasim: an agent-based model of COVID-19 dynamics and interventions. PLOS Comput. Biol. **17**(7), 1–32 (2021)
11. AnyLogic: Simulation Modeling Software Tools & Solutions for Business, https://www.anylogic.com/. Accessed 09 Feb 2023
12. Statistical Data of Ashiya, https://www.city.ashiya.lg.jp.e.akj.hp.transer.com/bunsho/toukei/index.html. Accessed 09 Feb 2023
13. Information for Medical Clinic (Treatment Guideline and Clinical Study) (in Japanese), https://www.mhlw.go.jp/stf/seisakunitsuite/bunya/0000121431_00111.html. Accessed 09 Feb 2023
14. National Urban Traffic Characteristics Survey Data (in Japanese), https://www.mlit.go.jp/toshi/tosiko/toshi_tosiko_fr_000024.html. Accessed 09 Feb 2023

# A Theoretical Agent-Based Model to Simulate the Rise of Complex Societies

**Saida Hachimi El Idrissi** [ID]**, Mohamed Nemiche** [ID]**, and Bezza Hafidi**

**Abstract** Nowadays, societies consist of hundreds of millions of people governed by one political system, and cooperation between individuals transcends face-to-face cooperation. However, in early history, groups of people did not exceed hundreds of individuals, and cooperation existed at low levels. So how did human societies evolve from small groups known by face and name into the huge anonymous groups of today? Our model tries to answer this question based on Freud's hypothesis stating that civilization could not arise and evolve without the repression of satisfaction (repression of human desires). In social evolution the repression of satisfaction can be interpreted as the repression of competition between society members the thing that increases the society power and helps on the formation of complex societies. In order to test this hypothesis we implemented an agent-based model where a large number of primitive societies are distributed in a grid of cells; initially, each cell is an independent polity. In each time step, all border cells (cells having at least one neighbor of a different polity) have a chance to start an attack and take over one cell from its neighbors. During the simulations we can observe the emergence of complex societies the thing that validates our hypothesis.

**Keywords** Agent-based modeling · Repression of satisfaction · Cooperation · Competition · Social simulation

## 1 Introduction

In the present age, the population of the largest countries is more than a billion. The world is full of societies with millions of citizens living in large lands ruled by one political system and maybe even by one person, the thing that was impossible in

---

S. Hachimi El Idrissi (✉) · B. Hafidi
IMI Laboratory, Ibn Zohr University, 80060 Agadir, Morocco
e-mail: saida.hachimi12@gmail.com

M. Nemiche
Polydisciplinary Faculty of Taza, 35000 Taza, Morocco

the early ages before agriculture. The emergence of agriculture allowed humans to settle in fertile lands and have a surplus of food, which facilitated the reproduction and growth of primitive groups. According to Turchin, the competition for resources between neighboring groups grew, as well as the attacks of nomadic groups aiming to get agricultural products by force. The competition between agriculturalists and nomads obliged the groups with the same interests to cooperate and defend their welfare. Consequently, the groups neglected their differences and built societies [1].

The evolution of all species is based on the natural selection [2]. However, doesn't natural selection favor competition and selfish behaviors more than cooperation? Then how could humans cooperate and organize huge societies where individuals have no genetic relations? Even with the advantage of selfishness and competition in nature, cooperation behavior has been observed, and many scientists tried to explain it. Cooperation exists at several levels. The first one is cooperation between relatives (kin selection), which is explained as the self-sacrifice for the favor of common genes [3]. The second is the cooperation between members in small groups (face-to-face cooperation), which in turn is explained by the reciprocity, reputation, retribution, and group selection [4–7].

The last level is large-scale cooperation, which is the key factor to the emergence of large-scale societies. Robert Boyed et al. argued that cultural adaptation is the main reason behind the evolution of large-scale cooperation [8]. They based their hypothesis on three assumptions. The first is that humans developed the ability to learn from each other, which created an evolution by cultural accumulation. Therefore, this ability was favored by natural selection because it provided the cultural adaptation needed in their social environment during the rapidly changing climates. The second is that rapid cultural adaptation increased the differences between groups and also increased competition between groups. However, reciprocity and reputation systems can balance selfish and cooperative behaviors within groups. The last assumption is that in culturally evolved cooperative groups, social selection and moral systems favor the reproduction of members with social norms that support pro-social motives; it also punishes the members who violate those norms, making their chances of reproduction very low. Thus, those moral systems forced by punishment and reward favored the success of individuals that functioned well in such an environment and also favored motives like shame, guilt, and other norms that facilitate the rise of large-scale cooperation [8].

In addition to cooperation, repression of competition also has an important role in the rise of large-scale societies. Alexander argued that repression of reproductive competition in human groups, added to the ferocity of competition between outside groups, helped spread human social structures [5, 9, 10]. Those social structures (social norms and social institutions) are the principal key to the emergence of large-scale societies. Frank goes with the same hypothesis and states that the fitness of a group increases with the decrease of competition within it. Low competition maximizes the number of individuals who benefit from resources and prevents the damage caused by overexploitation [10]. However, natural selection does not favor the repression of competition inside groups because individuals try to get maximum resources at the expense of their neighbors to ensure their survival. Besides that, we

can observe the evolution of internal repression traits in nature [10–14]. To explain that, Frank proposed a model where individuals have a variable called competition intensity (z) and another variable (a) called mutual policing, which represents the individuals' contribution to the repression of competition within the group. This mutual policing, also called punishment, helps the group reduce the competition and increase the cooperation among group members, and as a result, the group gains higher fitness [10].

The ability of humans to cooperate and live in huge groups without any genetic relations (this ability is called ultrasociality) produces large-scale societies [15]. A theoretical model presented by Turchin et al. [15] suggested that the evolution of ultrasocial norms and institutions is the reason behind the emergence of large-scale societies. As well, those ultrasocial norms and institutions result from intense competition and warfare between societies which depends on the spread intensity of military technology. Turchin used two vectors in his model: the first represents ultrasociality traits responsible for the power of a polity. The second vector is for military technology that is basically responsible for warfare intensity and ethnocide [15].

Another try to explain the huge cultural diversity in the world is the model of Talukdar et al. [16]. They made a computer model in a two-dimensional domain to simulate the interactions between cultures. The main dynamical processes used in their work are inspired by historical rules of expansion, interaction, and merging among cultures, namely growth, assimilation, invasion, aggression, and annihilation. From simple rules to define different interactions between cultures, the model pictured some interesting data in agreement with historical data, such as the intensity of wars in the primitive period compared to the modern one and the appearance of globally polarized cultures [16].

In this work, we propose a theoretical agent-based model to explain how large-scale societies emerged based on Freud's assumption about the emergence of civilization. Freud considers that without the restraint of human desires, civilization could not exist [17]. The central premise of this model is that the repression of individuals' satisfaction within the group that helps have strong societies is the result of the suffered repression caused by the neighboring societies (the risk of an attack from a neighboring group). The outside danger forces the individuals to cooperate and contribute to mutual policing to reduce competition inside the group.

Agent-based modeling (ABM) is a bottom-up approach of social simulation used to facilitate the modulization of social complex phenomenon. One of the main advantages of ABM is its ability to produce macro-scale phenomena (complex patterns) from simple behavioral rules at the micro-scale. It is based on agents who interact with each other and their environment according to simple rules. Agents may be individuals or collective entities that have attributes and methods. While the environment is the place where they live, they exploit it and interact with it according to specific rules. As for the rules, they are simple instructions that organize the interactions between the model elements [18–20].

## 2  Methods

### 2.1  General Logic of the Model

In this work, we develop an agent based model to understand how human societies grow from primitive societies where small groups are unified with face-to-face cooperation to complex societies of today. Our study is limited to the Old World, where the conflicts between societies occur only by land. Our simulation takes place on a two-dimensional hexagonal grid of cells. Each cell represents a community (local society). Each community is characterized by a satisfaction vector, a fitness, a power, and a suffered repression. At the start, each community is independent, and has six neighbors.

**Satisfaction** noted as $\boldsymbol{\pi^i}$, is a binary vector with $n_{sat}$ traits $\pi^i(t) = \left(\pi_1^i(t), \ldots, \pi_{n_{sat}}^i(t)\right)$, where $\pi_k^i(t) \in \{0, 1\}, \forall k \in \{0, \ldots, n_{sat}\}$, where $n_{sat}$ is a parameter of the model. The presence of all the satisfaction traits in a polity refers to a primitive society where there is no norms or restrictions to respect.

**The satisfaction intensity** for an independent polity "i" is calculated as the average value of the satisfaction traits:

$$\overline{\pi}^i(t) = \frac{1}{n_{sat}} \sum_{l=1}^{n_{sat}} \pi_l^i(t) \tag{1}$$

For primitive society $\overline{\pi}^i(t) = 1$.

For multicell polity "i" the satisfaction intensity is defined as:

$$\overline{\Pi}_i(t) = \frac{1}{S_i} \sum_{k=1}^{S_i} \overline{\pi}^k(t) \tag{2}$$

where $\boldsymbol{S_i}$ is the polity size (number of cells of the multicell polity "i"), and $\overline{\pi}^k$ the satisfaction intensity of the individual polity "k" [21].

To calculate the value of the individual fitness of a polity "j" that belongs to a multicell polity "i", we are inspired by Frank's formula [10]:

$$w_{ij}(t) = \left(\overline{a}_i(t) - c \cdot a_{ij}(t) + (1 - \overline{a}_i(t)) \cdot \frac{\overline{\pi}^j(t)}{\overline{\Pi}_i(t)}\right) \left(1 - (1 - \overline{a}_i(t)) \cdot \overline{\Pi}_i(t)\right)(1 - \overline{\sigma}_i(t)) \tag{3}$$

where

- $\overline{\sigma}_i(t)$ represents the suffered repression by the multicell polity "i" from its neighbors (social context) [12–14, 22];

- $a_{ij}(t)$ is an individual's participation in mutual policing. Policing is a mechanism that reduces the competition within the group by repressing the satisfaction of individuals. The value of the individual's policing increases with the increase of the suffered repression (external danger); which favors intragroup cooperation and increases the fitness of the group [10];
- $\overline{a}_i(t)$ is the average level of policing in the polity "i";
- $ca_{ij}(t)$ is the cost to live in a group;
- c is a parameter of the model;

For one-cell polities, the fitness is simplified as:

$$w_i(t) = \left(1 - \overline{\pi}^i\right)(1 - \overline{\sigma}_i(t)) \tag{4}$$

According to Freud, repression of individual satisfaction is necessary for technical progress of the group [17]. Which leads to powerful polities and facilitates the rise of complex societies.

**The individual policing** increases with the increase of suffered repression and the increase of the average satisfaction:

$$a_{ij}(t) = \overline{\pi}^j(t).\overline{\sigma}_i(t) \tag{5}$$

**The Suffered Repression** $\overline{\sigma}_i(t)$ represents the danger of neighboring societies [12–14, 22]. It is calculated based on their Power:

$$\overline{\sigma}_i(t) = \frac{\sum_{j \in V_i^*}^{j} Power_j(t-1)}{\sum_{j \in V_i}^{j} Power_j(t-1)} \tag{6}$$

$V_i^*$ is the set of neighboring societies of the polity "i" with "i" excluded, $V_i$ is the same neighborhood with polity "i" included. At the start of the model we consider $\overline{\sigma}_i(0) = 0$ for all cells.

**Power** of a polity "i" is defined as [15]:

$$Power_i(t) = 1 + \beta S_i \overline{w}_i(t) \tag{7}$$

with

$$\overline{w}_i(t) = \frac{1}{S_i} \sum_{j=0}^{S_i} w_{ij}(t) \tag{8}$$

both the size $S_i$ and the average fitness of the polity $\overline{w}_i$ increase the polity's Power, $\beta$ is the coefficient that translates fitness into polity's power.

We summarize the variables of our model in Table 1:

**Table 1**  Entities of our model

| Entities | Variable name | Possible values |
|---|---|---|
| One cell polity "i", it can be a piece of an empire"j", Or An independent polity (local polity "i") | (x, y) coordinates (localization) | ([− 27, 27], [− 27, 27]) |
| | Imperial-index | 0: for independent polity<br>An integer for polities that belong to an empire {1, … 100,000} |
| | Satisfaction vector $\boldsymbol{\pi^i}$ | $\pi^i(t) = \left(\pi^i_1(t),\ \ldots,\ \pi^i_{n_{sat}}(t)\right)$,<br>where: $\pi^i_k(t) \in \{0, 1\} \forall k \in \{0, \ldots, n_{sat}\}$ |
| | Satisfaction intensity $\boldsymbol{\overline{\pi}^i}$ (average value of satisfaction traits) | [0.1,1] |
| | Individual fitness $\boldsymbol{w_{ij}}$ Or ($\boldsymbol{w_i}$ for independent polity) | [0, 1] |
| | Power: $Power_i$ calculated for independent polity | [1, 500] |
| | Suffered repression: $\overline{\sigma}_i$ calculated for independent polity | [0, 1] |
| | Individual policing $\boldsymbol{a_{ij}}$, ( null for independent polity) | [0, 1] |
| Multicell polity (Empire "i") | Imperial-Index | {1, …, 100,000} |
| | Average satisfaction intensity: $\overline{\Pi}_i$ | [0.1, 1] |
| | Average fitness: $\overline{\boldsymbol{w}}_i$ | [0, 1] |
| | Power: $\boldsymbol{Power_i}$ | [1, 500] |
| | Suffered repression $\overline{\sigma}_i$ | [0, 1] |
| | Average policing $\overline{\boldsymbol{a}}_i$ | [0, 1] |
| | Size $\boldsymbol{S_i}$ | {1, 2, …, 600} |
| Environment | Grid of hexagonal cells | 55 × 55 = 3025 cells |

## 2.2  Warfare Between Polities

In this model, conflicts between societies are managed in the same way as in Turchin's model [15]. In each time step, all border cells (cells having at least one neighbor of a different polity) have a chance to start an attack in a random direction. If the attacked cell is from the same polity nothing happens, if not, war can be initiated between two polities with probability P-attack, with P-attack a parameter of the model. The order in which the attacker cells are chosen is randomized every time step [15].

An attack can be successful with probability $P_{success}$ defined as:

$$P_{success}(t) = \frac{Power_{attacker}(t) - Power_{defender}(t)}{Power_{attacker}(t) + Power_{defender}(t)} \tag{9}$$

If $P_{success} < 0$ the attack fails by definition, and nothing happens. If $P_{success} > 0$ the attack is successful, the attacked cell will be annexed to the polity of the attacker. The attacked cell may also copy the satisfaction vector of the attacker, with a probability $P_{ethnocide}$.

## 2.3 Sociocultural Dynamics: (Mutation, Ethnocide)

The dynamic process of our satisfaction vector is defined by two mechanisms:

- The first one is mutation, it represents random changes in the satisfaction vector [15]. At each time step, for every cell, we chose a random position in satisfaction vector; its value may change from 0 to 1 with a probability $\mu_{01}$ if the value of the trait is 0, or from 1 to 0 with probability $\mu_{10}$ if the trait value is 1. We note that we made the first position of all satisfaction vectors equal 1, and we don't change them in mutation, so we eliminate the case where $\overline{\Pi}_i(t) = 0$. We assume that gaining a satisfaction trait is much easier than losing it $\mu_{01} \gg \mu_{10}$ because it is easier to break social norms than to respect them and add new ones [15].
- The second mechanism is Ethnocide or (Forced cultural assimilation). After an attack, we calculate a probability [15]:

$$P_{ethnocide}(t) = \varepsilon_{min} + (\varepsilon_{max} - \varepsilon_{min})\left(1 - \overline{\Pi}_{attacker}(t)\right) \tag{10}$$

and the defeated cell may copy the satisfaction vector of the winning cell with probability $P_{ethnocide}$. $\varepsilon_{min}$ is the minimum value $P_{ethnocide}$ can have when the average satisfaction is at maximum 1, and $\varepsilon_{max}$ is the maximum value when the average satisfaction is near to 0. Here we consider that a society with less satisfaction can have more control over new annexed cells.

## 2.4 Collapse

Wars cause societies to grow and expand through assimilation, aggression and annihilation. However, they all inevitably decay afterwards because of the repercussions of those wars and the pluralism leading to civil wars, among other factors. In our model we define the probability of collapse as in Turchin model [15]. At each time step, each polity can disintegrate into polities with one cell territory each. The probability of the collapse increases with the polity size $S_i$ and decreases with the average fitness $\overline{w}_i$:

$$p_i(t) = \delta_0 + \delta_1 S_i - \delta_2 \overline{w}_i(t) \qquad (11)$$

$\delta_0$, $\delta_1$ and $\delta_2$ are parameters of the model. $\delta_0 = 1/20$ represents the baseline disintegration probability. $\delta_1 = 1/20$, so a polity with low fitness will certainly collapse after size of 19 cells. And finally the $\delta_2$ is defined in a way that every empire that reaches the size 600 will collapse [15].

## 3   Results and Discussion

The simulation of our model shows, in a two-dimensional grid of cells, the formation, the expansion, and the collapse of many complex societies. The main result of this work is as expected: the decrease of individual satisfaction facilitates the rise of large-scale societies.

From the simulations, we observed that societies appear in the same way across all the chronology of our simulations, where we note the emergence, expansion, and then decline of complex societies. However, the difference lies in the size of the constituent empires as they expand more as time progresses in the model. At the start of the simulation, the formed societies are small with relatively high satisfaction (which translates the absolute individual freedom in primitive societies where there are no norms and laws to organize communal life), afterward as time passes, due to the logic of our model, the societies' satisfaction decreases the thing that refers to the emergence of social norms and institutions that enable human groups to cooperate with each other and live in huge societies.

Here we present images from a simulation as an example to show the general pattern of the simulations.

Every color defines a different society, and the size is the number of cells with the same color.

To test our model, we decide to examine how sensitive our results are according to the variations of the variables: Number of Traits "$n_{sat}$" of the satisfaction vector, and the cost to live in an Empire "c". For other parameters, we chose fixed values as a start to facilitate our tests.

- Note: for figure Fig. 1 (images a, b and c) the values used for "$n_{sat}$" and "c" are respectively 10 and 0.6.

Since our essential purpose in this model is the emergence of complex societies, we focused our interest on societies with a size of more than one hundred cells (which we can consider as empires). To facilitate tracking those empires, we divided the timeline (1500 time step) into periods of 100-time step each. At the end of the simulation, we got all the empires (societies with 100 cells or more) formed for each period with their necessary data. From that, we can compare between different periods and determine the variation of empires through time.

The values we chose for our parameters are: $n_{sat} = \{4; 7; 10\}$, $c = \{0.3; 0.6; 0.9\}$.

A : After 40 time steps    B : After 360 time steps    C : After 700 time steps

**Fig. 1** **a** Image of a simulation after 40 time steps, we can only see small societies. **b** Image of a simulation after 360 time steps, we can see the rise of big societies. **c** In this image we can observe the dominance of two big empires at the time step 700

After varying values of $n_{sat}$ and $c$, we got a set of simulation results where we concluded that our model kept the same logic for all the parameter combinations (we can still observe the emergence, expansion, and collapse of complex societies).The difference is in the average sizes of the constructed empires and the ability to decrease the satisfaction intensity of societies (the thing that facilitates the rise of the size).

We observe from the results that higher values of $n_{sat}$ decrease the values of satisfaction intensity (see Fig. 4), and give empires with large sizes (see Fig. 2). However, the variable $c$ has no direct impact on satisfaction and size variations because, after many periods, their respective values stabilize at the same value each for all the variations of c (see Figs. 3 and 5). On the other hand, the variable $c$ affects the number of empires observed in each period so that their number decreases with the increase of the variable $c$ (Figs. 6 and 7).

- Note: For all the graphs the number of simulations is 30, each one has 1500 unit of time. Each period equals 100 unit of time.



**Fig. 2** Variation of the average size of empires for each period of time and for each value of the variable "$n_{sat}$"—for each value of $n_{sat} = \{4, 7, 10\}$ we calculate the average values of all the variations of c = 0.3, 0.6 and 0.9

**Fig. 3** Variation of the average size of empires for each period of time and for each value of the variable "c"—for each value of c = {0.3, 0.6, 0.9} we calculate the average values of all the variations of $n_{sat}$ = 4, 7 and 10



**Fig. 4** Variation of the satisfaction intensity of empires for each period of time and for each value of the variable "$n_{sat}$"—for each value of $n_{sat}$ = {4, 7, 10} we calculate the average values of all the variations of c = 0.3, 0.6 and 0.9



**Fig. 5** Variation of the satisfaction intensity of empires for each period of time and for each value of the variable "c"—for each value of c = {0.3, 0.6, 0.9} we calculate the average values of all the variations of $n_{sat}$ = 4, 7 and 10

**Fig. 6** Variation of the Number of Empires that emerged for each period of time and for each value of the variable "$n_{sat}$"—for each value of $n_{sat} = \{4, 7, 10\}$ we calculate the average values of all the variations of $c = 0.3, 0.6$ and $0.9$



**Fig. 7** Variation of the Number of Empires that emerged for each period of time and for each value of the variable "$c$"—for each value of $c = \{0.3, 0.6, 0.9\}$ we calculate the average values of all the variations of $n_{sat} = 4, 7$ and $10$

## 4 Conclusion

In this paper, we present our work consisting of a simple agent-based model that explains how large-scale societies emerged in the old world. We have considered war, and repression of satisfaction as important mechanisms that control cooperation within groups. Warfare and competition between societies enhance cooperation inside groups aiming to defend their territories and common interests. Thus, the evolution of cooperation inside groups facilitates the group's expansion at the expense of other groups and hence promotes the rise of large societies.

Our work could be extended by adding geographical parameters and reimplementing the model in a realistic geographical environment. These geographical parameters will specify the type of territories: fertile lands, mountains, rivers, and seas, and will define places with intense warfare and places where the defense is easier compared to others.

# References

1. Turchin, P.: A Theory for the Formation of Large Agrarian Empires. Santa Fe Institute Working Paper, p. 05 (2008)
2. Darwin, C.: On the Origin of Species: A Facsimile of the First Edition. Harvard University Press (1964)
3. Hamilton, W.D.: The genetical evolution theory of social behaviour, I and II. J. Theor. Biol. **7**(1), 1–52 (1964)
4. Trivers, R.L.: The evolution of reciprocal altruism. Q. Rev. Biol. **46**(1), 35–57 (1971)
5. Alexander, R.D.: The Biology of Moral Systems. Routledge (2017)
6. Haley, K.J., Fessler, D.M.: Nobody's watching? Subtle cues affect generosity in an anonymous economic game. Evol. Hum. Behav. **26**(3), 245–256 (2005)
7. Nowak, M.A., Sigmund, K.: Evolution of indirect reciprocity. Nature **437**(7063), 1291–1298 (2005)
8. Boyd, R., Richerson, P.J.: Culture and the evolution of human cooperation. Philos. Trans. R. Soc. B Biol. Sci. **364**(1533), 3281–3288 (2009)
9. Alexander, R.D.: Darwinism and Human Affairs. University of Washington Press, Seattle (1979)
10. Frank, S.A.: Mutual policing and repression of competition in the evolution of cooperative groups. Nature **377**(6549), 520–522 (1995)
11. Szathmáry, E., Smith, J.M.: The major evolutionary transitions. Nature **374**(6519), 227–232 (1995)
12. Nemiche, M.: Un Modelo Sistémico de Evolución Social Dual. Doctoral thesis, Universidad de Valencia, Spain (2002)
13. M'hamdi, A., Sfa, F.E., Nemiche, M., Hachimi El Idrissi, S., Pla-López, R.: Modelling "Occident/Orient" duality and migration process with mobile agents. Syst. Res. Behav. Sci. 36(6), 750–764 (2019)
14. M'hamdi, A., Nemiche, M.: Bottom-up and top-down approaches to simulate complex social phenomena. Int. J. Appl. Evol. Comput. 9(2), 1–16 (2018)
15. Turchin, P., Currie, T.E., Turner, E.A., Gavrilets, S.: War, space, and the evolution of old world complex societies. Proc. Natl. Acad. Sci. **110**(41), 16384–16389 (2013)
16. Talukdar, D., Dutta, K.: An archetype for evolving dynamics of primitive human culture. Evol. Syst. **12**(4), 965–979 (2021)
17. Freud, S.: Civilization and Its Discontents. Broadview Press (2015)
18. Epstein, J.M., Axtell, R.: Growing Artificial Societies: Social Science from the Bottom Up. Brookings Institution Press (1996)
19. Cioffi-Revilla, C.: Introduction to Computational Social Science. Springer, London and Heidelberg (2014)
20. M'hamdi, A., Nemiche, M., Pla Lopez, R., Ezzahra SFA, F., Sidati, K., Baz, O.: A generic agent-based model of historical social behaviors change. In: Nemiche, M., Essaaidi, M. (eds.) Advances in Complex Societal, Environmental and Engineered Systems, pp. 31–49. Springer, Cham (2017)
21. El Idrissi, S.H., Nemiche, M., Chakraoui, M.: Repression of satisfaction as the basis of the emergence of old world complex societies. Compl. Syst. **29**(3), 655–667 (2020)
22. Nemiche, M., M'Hamdi, A., Chakraoui, M., Cavero, V., Pla Lopez, R.: A theoretical agent-based model to simulate an artificial social evolution. Syst. Res. Behav. Sci. **30**(6), 693–702 (2013)

# An Agent-Based Model of Prosocial Equilibrium: The Role of Religiously Motivated Behaviour in the Formation and Maintenance of Large-Scale Societies

**Ivan Puga-Gonzalez** 🅞**, F. LeRon Shults** 🅞**, Ross Gore** 🅞**, and Konrad Talmont-Kaminski** 🅞

**Abstract** This paper outlines a new agent-based model that represents the dynamics by which a society achieves and maintains prosocial equilibrium. The latter is understood as a social balance involving the interplay of prosocial behavior, anxiety, environmental threats, and religiosity in the population. Experiments showed that the model was able to simulate the emergence of relatively large societies under the sorts of conditions that would be expected based on the theoretical literature and other empirical findings in the relevant fields. We conclude by describing the main insights of the simulation experiments and pointing toward future work currently being planned by the research team.

**Keywords** Agent-based model · Religiosity · Prosocial behavior

## 1 Introduction

### 1.1 Theory on Prosocial Behavior and Intro to the Model

The study of the behavior of humans and other animals in the last few decades has spent considerable efforts looking at explanations of altruistic behavior. Kin selection and reciprocal altruism have helped to understand cooperation among other animals

I. Puga-Gonzalez (✉) · F. L. Shults
Center for Modeling Social Systems, NORCE, Kristiansand, Norway
e-mail: ivanpuga@gmail.com

R. Gore
Virginia Modeling, Analysis and Simulation Center, Old Dominion University, Norfolk, VA, USA

K. Talmont-Kaminski
Society and Cognition Unit, University of Bialystok, Białystok, Poland

F. L. Shults
Institute for Global Development and Social Planning, University of Agder, Kristiansand, Norway

[1, 2] and in small-scale human societies, without reference to previously popular but largely discredited group selection explanations [3]. These results have, however, thrown into sharp contrast the situation with large-scale human societies. The issue is that in large-scale societies, where most interactions are one-off with unrelated individuals, the logic behind the theories used to explain animal cooperation breaks down and the strategy of becoming a free rider appears to be much more attractive [4]. Many solutions were suggested for this problem, but most of these were largely concerned with altruistic punishment [5], which was odd because on the whole this was still a form of altruistic behavior that was costly to the individual [6].

The lack of a generally accepted explanation for cooperation in large-scale human societies motivated some researchers to look to religion for an explanation. The idea that religion could play the role of motivating social cohesion and cooperation is one that has a long history, with numerous scholars having argued for it [7, 8]. In recent years, scholars working within cognitive science of religion have proposed several mechanisms by which religious beliefs and practices could help to motivate costly pro-social behavior [9, 10] as well as seeking to connect the historical appearance of large-scale societies with changes in religious traditions that could have played a role in making those societies possible [8, 11]. At the same time, a range of studies also going back close to a hundred years has provided evidence for the claim that increased levels of anxiety lead—in the short-term as well as well in the long-term—to increased espousal of religious claims and engagement in religious practices [12–14].

When taken together, the connection between religion and cooperation as well as the connection between anxiety and religion potentially form two parts of a negative feedback mechanism that could underpin a prosocial equilibrium, the connection between cooperation and anxiety forming the final element [15, 16]. The picture is that of environmental threats leading to increased anxiety and thereby to increased religiosity. However, an increase in religiosity drives increased cooperation, which helps the society deal with the threats and thereby lower anxiety. In effect, a relatively high level of religiosity and cooperation is maintained, allowing large-scale societies to thrive.

To test the plausibility of such a prosocial equilibrium we decided to construct an agent-based model in which altruistic, prosocial behavior that was not motivated by reciprocal arrangements or genetic connections would have the opportunity to allow the growth of societies with many hundreds or even thousands of members. A key assumption of the model was that religiosity is primarily determined by observing/participating in religious practices during the period of socialization. The practices focused upon in the model are forms of religiously-motivated pro-social behavior such as participation in work for the community, tithing and other forms of religious charity, as well as those forms of sacrifice in which the offering is made use of by the community—which are considered to play a large role in motivating religiosity.

## 2   Methods

### 2.1   The Model

The model was written in AnyLogic v.8.7.9. Here we present a brief description of the model. A full ODD + D protocol description can be found at the github repository: URL https://github.com/ivanpugagonzalez/Prosociality-ABM-Model

**Model overview and agents**. The artificial society represented in the model is inhabited by individual human agents who have eight different variables: age, gender, marital status, religiosity, wellbeing, insecurity, sensitivity, and anxiety. On initialization, 1000 adult agents are created, the age distribution follows a typical pyramid shape (0–100 years). The initial values of religiosity, and sensitivity are drawn from a normal distribution $N(\mu = 0.5, \sigma = 0.1)$ and that of insecurity is set to 0. Every year all agents experience a given number of environmental threats of different intensity. The number and intensity of threats are controlled by a Poisson and exponential distribution (13–14 in Table 1). Threats increase the insecurity of agents and in turn insecurity increases anxiety. Anxiety and religiosity may then trigger a prosocial behavior. Prosocial behaviors increase the religiosity and decrease the insecurity of the performing agents and that of close by neighbors (7–10 in Table 1). Prosocial behaviors are costly and reduce the wellbeing of performing agents (11 in Table 1). Agents also increase/decrease their wellbeing according to their current age and insecurity values (15–21 in Table 1). Agents reproduce if they are married, female, and within the age of reproduction. Agents that are 25 y.o. or younger, reduce a percentage of their religiosity every year (24 in Table 1). Agents die with a probability given by their wellbeing value. Figure 1 shows a summary of the model cycle and order of processes during the simulation.

**Wellbeing processes**. Wellbeing (WB) determines the probability of an agent surviving every year. A survival probability curve was mimic using data from 1950's in Norway. This choice was arbitrary, but it doesn't have a major effect on the model's behavior. Both the reference model and the one with prosocial behavior (see Sects. 2.2 and 2.3) use the same survival probability curve, and because we compare one against the other the effect of the survival probability curve becomes irrelevant.

*Wellbeing and age.* At initialization, wellbeing is determined by a polynomial function of the agents' age. This equation mimics the survival probability of both sexes according to age during 1950's in Norway. Then, after initialization, WB of agents increases and decrease every year according to its age. The gain or loss of WB is determined by two equations.

The gain of WB is given by equation 1:

$$Gain = -4C * \left( \frac{Age - WB\_Age\_Threshold}{100 - WB\_Age\_Threshold} \right)^{Exp1} + C \qquad (1)$$

**Table 1** Model parameters

| Parameter | Value | Description | Process |
|---|---|---|---|
| 1. Rep Cost | OP | % of WB taken from each parent | Rep |
| 2. Rep mid threshold | OP | Reproduction probability is 0.5 | |
| 3. Rep Curve Shape | OP | Parameters determining the shape of probability of reproduction curve | |
| 4. Importance Insec | SA | | |
| 5. Importance WB | 1 | | |
| 6. PB threshold | SA | Threshold value to trigger PB | PB |
| 7. PB inc rel self | SA | Increase in agent's and neighbors' religiosity after a PB | |
| 8. PB inc rel neigh | SA | | |
| 9. PB dec insec self | SA | Decrease in agent's and neighbors' insecurity after a PB | |
| 10. PB dec insec neigh | SA | | |
| 11. PB wellbeing cost | SA | Decrease in agent's WB after a PB | |
| 12. Neigh Benefited | SA | # of nearby neighbors benefited | |
| 13. Threats Rate | SA | Shape of the Poisson distribution | Threats |
| 14. Threats Intensity | SA | Shape of the exponential distribution | |
| 15. WB Age Threshold | OP | Parameters determining the increase/ decrease of WB according to agents' age | WB-Age |
| 16. WB Intercept C | OP | | |
| 17. WB Exp Gain eq | OP | | |
| 18. WB Exp Loss eq | OP | | |
| 19. WB Insec Threshold | 0.1 | Parameters determining the increase/ decrease of WB according to agents' insecurity | WB-Insecurity |
| 20. WB Max Inc | OP | | |
| 21. WB Max Dec | 0.25 | | |
| 22. Marriage Age Diff | OP | Max age difference between partners | Others |
| 23. Radius Local Area | 50 | Radius of area of nearby neighbors | |
| 24. Rel Dec Perc | SA | % of religiosity decrease every year | |

*WB* wellbeing, *PB* prosocial behavior, *Insec* insecurity, *Rep* reproduction, *inc* increase, *dec* decrease, *rel* religiosity, *OP* optimized parameter, *SA* sensitivity analysis

The loss in WB is then given by equation 2:

$$Loss = -4C * \left( \frac{Age - WB\_Age\_Threshold}{100 - WB\_Age\_Threshold} \right)^{Exp2} + C \qquad (2)$$

where WB Age Threshold (15 in Table 1) is the age at which the gain/loss in WB is given by equation 2 instead of equation 1; C (16 in Table 1) is the equation intercept, and *Exp1* and *Exp2* (17–18 in Table 1) determine the shape of the curve.

*Wellbeing and insecurity.* In addition to being affected by age, WB is also affected by the agents' insecurity. Depending on the insecurity value and the value of *WB Insec Threshold* (19 in Table 1), wellbeing may increase or decrease every year according to equations 3 and 4 respectively.

**Fig. 1** Model cycle and order of processes

If insecurity ≤ *WB Insec Threshold*:

$$Gain = WB.Max.Inc + \left( Ins * \frac{WB.Max.Inc}{WB.Insec.Th} \right) \qquad (3)$$

The left-hand side term corresponds to the intercept and the fraction on the right-hand side to the slope. *Ins* is the current insecurity of the agent; *WB.Max.Inc* represents the maximum gain in WB when insecurity equals 0; and *WB.Ins.Th* is the insecurity value at which there is neither gain nor loss in WB (19–20 in Table 1).

If insecurity > *WB Insec Threshold:*

$$Loss = \frac{-WB.Max.Dec}{(1 - WB.Insec.Th)} * WB.Insec.Th + \left( Ins * \frac{WB.Max.Dec}{(1 - WB.Insec.Th)} \right)$$

(4)

The left-hand side term corresponds to the intercept and the fraction on the right-hand side to the slope. *Ins* is the current insecurity of the agent; *WB.Max.Dec* represents the maximum loss in WB when insecurity equals 1; and *WB.Ins.Th* is the insecurity value at which there is neither gain nor loss in WB (19–20 in Table 1).

**Mortality process**. As previously mentioned, WB determines the probability of agents dying and mimics the probability of dying reported in census data during 1950's in Norway. According to this data, the probability of dying for both sexes increase with age (Fig. 2a). To mimic this probability, we fitted a polynomial curve across the census data and input wellbeing instead of age. This resulted in the probability of dying curve shown in Fig. 2b.

**Marriage and Reproduction process**. To marry, agents had to meet several conditions: not being married, being over 15 y.o., and that the age difference between potential partners is not higher than *Marriage Age Diff* (22 in Table 1). If these conditions were met, agents were set to a married marital status.

Once married, female agents in the age of reproduction have the chance to reproduce every year. The probability of reproduction depends on the WB and insecurity of the married agents, it is given by equation 5:

$$Prob.Rep = \frac{1}{1 + e^{(-b*(x-a))}}$$

(5)



**Fig. 2** Probability of dying according to **a** census data and **b** wellbeing

where *b* is the parameter *Rep curve shape* (3 in Table 1) determining the shape of the sigmoidal curve. *x* is a weighted average equal to:

$$\frac{(Average.WB.Partners) * Importance.WB + (Average.Ins.Partners) * Importance.Ins}{Importance.WB + Importance.Ins}$$

and represents the importance of WB and insecurity in the reproduction decision (4–5 in Table 1), and *a* is the WB threshold at which reproduction probability is equal to 0.5 (2 in Table 1).

If agents reproduce, then their WB is decreased by a percentage given by Rep Cost (1 in Table 1). The loss in WB from both partners is then passed into the offspring, and this value becomes the initial WB value of the offspring. Further, offspring inherit the religiosity, insecurity, and sensitivity values from one of their parents (this parent is selected at random).

**Threats process**. Every year a certain number of threats are generated. The number and intensity of threats is determined by drawing numbers from a Poisson distribution and an exponential distribution, respectively (13–14 in Table 1). Each year, the intensity value of all threats is added to the current insecurity of agents (if after this addition insecurity is > 1, insecurity is set to 1).

**Prosocial behavior process**. Every year, agents aged 12 y.o. or older are allowed to perform a prosocial behavior. Prosocial behavior is triggered when the product of the agents' religiosity times anxiety goes above the *PB threshold* (8 in Table 1). Anxiety is the product of their insecurity times their sensitivity. Prosocial behaviors increase the religiosity and decrease the insecurity of the performing agents and that of close by neighbors (7–10 in Table 1). Neighbors are considered those agents within a certain radius of distance from the performing agent (24 in Table 1). If the number of close by agents exceed the number of benefited neighbors (23 in Table 1); then, benefited neighbors are selected randomly from all close by ones. Prosocial behaviors are costly and performing agents reduce their wellbeing every time they perform a prosocial behavior (11 in Table 1).

## 2.2 Optimization Experiments and Reference Models

We created a reference model (RM) against which we could compare the effects of environmental threats and prosocial behavior on the growth rate of the society. The RM was needed because otherwise we would not know if societies were not successful because rate of threats and their intensity did not trigger the appropriate amount of prosocial behavior or because the parameters determining the mortality, reproduction and marriage processes were not tuned accurately and thus made societies go extinct. To create a RM, we turned off the environmental threats and prosocial

behaviors, and searched for optimal parameter values (related to wellbeing, mortality, marriage, and reproduction; OP parameters in Table 1) that allowed a society to keep its population size constant over time.

We used the optimization engine of AnyLogic, which allows the user to obtain a combination of parameter values that increases or decreases a specific output value obtained from an input function. In our case, the input function calculated the residual sum of squares (RSS) between the observed yearly growth rate (i.e., $pop\_size_{y+1}/pop\_size_y$) and the expected growth rate if the population size remained constant over time (i.e., 1). The optimization experiments found the combination of parameter values that minimize the output value. We ran 20 different optimization experiments from which we obtained 20 different combinations of parameters. Each simulation lasted for 500 years. We chose the two best simulations as RM (see ODD + D protocol for further details on the RMs).

## 2.3 Sensitivity Analysis: The Role of Threats and Prosocial Behavior

To explore the effect of threats and threats' intensity on prosocial behavior, religiosity, and the growth rate of societies, we did a sensitivity analysis by varying the values of 11 parameters related to prosociality, religiosity, reproduction, and threats (SA parameters in Table 1). During the sensitivity analysis we kept fixed the optimized parameters found for the two best simulations during the optimization experiments.

We used latinhypercube sampling to sample the parameter space 20,000 times per RM. For each combination of parameter values we ran one simulation. Each simulation was run for 500 years and every 25 years we collected the population size and average religiosity of the population. We classify as successful societies those that at the end of the simulation (500 years), had population size greater than 2000 individuals. We choose this value because it is a value greater than the median and the third interquartile range of population sizes of the RMs.

## 3 Results

### 3.1 Successful Societies

The percentage of simulations with successful societies (i.e., with a population size > 2000) for RM 1 and 2 were 0.44% (n = 88) and 2.15% (n = 431) respectively. We used the verification and validation (V&V) calculator tool (available at https://vmasc.shinyapps.io/VandVCalculator/), whose use is illustrated in [17–19], in order to explore the conditions leading to successful societies in our model. The Sensitivity Assessor identified four conditions that were observed much more frequently in the

successful societies than in all the parameter sampling. These conditions were related to specific parameters' values being below or above a certain threshold (Table 2). For instance, the PB threshold below 0.2 (or 0.3) was observed in 100% of the successful runs in RM1 (or 97% in RM2); whereas this condition was observed only in 39% of cases in the whole parameter sampling for RM1 (or 59% for RM2; Table 2). Similarly, the four conditions identified by the Sensitivity Assessor were observed over 90% of the time in successful simulations, a percentage well above the value expected from their appearance in the parameter sampling (Table 2). This suggested that societies may be successful if the values of PB threshold, PB wellbeing cost and the importance of insecurity in reproduction were kept below these thresholds and the number of neighbors benefited were above or equal to 2. Note that although these thresholds were somewhat different depending on the RM, the same parameters were identified in both models (Table 2).

Additionally, the Sensitivity Assessor identified other conditions concerning the value of the PB threshold and the WB cost of PB in relation to the value of other parameters (Table 3). More specifically, it seemed like successful societies were those where the values of both the PB threshold and WB cost of PB were lower than the decrease of insecurity (in self and neighbors) and the increase in religiosity (in self and neighbors) after a PB (Table 3).

To corroborate that these conditions were necessary for societies to be successful, we resampled the parameter. We first resampled (20,000 times per reference model) the parameter space by keeping the values of the parameters in Table 2 within the range identified by the Sensitivity Assessor. This indeed increased the percentage of successful societies: 12.59% (n = 2517) and 28.27% (n = 5655) for RM 1 and 2, respectively. This was a significant increase; however, the percentage of successful societies was far from being a majority, i.e., > 50%. Therefore, we decided to resample

**Table 2** Conditions observed in successful societies

| Condition | # of times observed in: | | % Obs | % Exp | Obs–Exp |
|---|---|---|---|---|---|
| | Successful runs | All sampling | | | |
| *Reference model 1* | | | | | |
| PB threshold < 0.2 | 88 | 7755 | 100 | 39 | 61 |
| PB wellbeing cost < 0.12 | 84 | 11,960 | 95 | 60 | 36 |
| Importance Insec < 0.8 | 86 | 2500 | 98 | 13 | 85 |
| # Neigh benefited > = 2 | 88 | 2500 | 100 | 88 | 13 |
| *Reference model 2* | | | | | |
| PB threshold < 0.3 | 417 | 11,836 | 97 | 59 | 38 |
| PB wellbeing cost < 0.14 | 395 | 13,969 | 92 | 70 | 22 |
| Importance Insec < 1.0 | 429 | 3333 | 100 | 17 | 83 |
| # Neigh benefited > = 2 | 416 | 2500 | 97 | 88 | 09 |

**Table 3** Extra conditions observed in successful societies

| Condition | # of times observed in: | | % Obs | % Exp | Obs–Exp |
|---|---|---|---|---|---|
| | Successful runs | All sampling | | | |
| *Reference model 1* | | | | | |
| PB th < PB dec ins neigh | 85 | 10,046 | 0.97 | 0.50 | 0.47 |
| PB th < PB dec ins self | 83 | 9999 | 0.94 | 0.50 | 0.44 |
| PB th < PB inc rel self | 85 | 10,025 | 0.97 | 0.50 | 0.47 |
| PB th < PB inc rel neigh | 77 | 10,065 | 0.88 | 0.50 | 0.38 |
| PB WB cost < PB dec ins neigh | 82 | 10,014 | 0.93 | 0.50 | 0.43 |
| PB WB cost < PB dec ins self | 74 | 10,001 | 0.84 | 0.50 | 0.34 |
| PB WB cost < PB inc rel self | 81 | 10,000 | 0.92 | 0.50 | 0.42 |
| PB WB cost < PB inc rel neigh | 72 | 10,007 | 0.82 | 0.50 | 0.32 |
| *Reference model 2* | | | | | |
| PB th < PB dec ins neigh | 330 | 10,014 | 0.77 | 0.50 | 0.27 |
| PB th < PB dec ins self | 346 | 9964 | 0.80 | 0.50 | 0.30 |
| PB th < PB inc rel self | 354 | 10,027 | 0.82 | 0.50 | 0.32 |
| PB th < PB inc rel neigh | 334 | 10,081 | 0.77 | 0.50 | 0.27 |
| PB WB cost < PB dec ins neigh | 316 | 10,010 | 0.73 | 0.50 | 0.23 |
| PB WB cost < PB dec ins self | 329 | 10,004 | 0.76 | 0.50 | 0.26 |
| PB WB cost < PB inc rel self | 329 | 10,005 | 0.76 | 0.50 | 0.26 |
| PB WB cost < PB inc rel neigh | 301 | 9973 | 0.70 | 0.50 | 0.20 |

the parameter space by not only maintaining the conditions in Table 2 but also those in Table 3. This resulted in the largest number of successful societies: 63.4% (n = 12,680) and 71.81% (n = 14,361) for RM 1 and 2 respectively.

## 3.2 *The Effect of Threats' Rate and Intensity*

Surprisingly, the conditions identified by the Sensitivity Assessor (Tables 2 and 3) did not include the rate and intensity of threats. However, we know that the rate and intensity of threat must play a role, otherwise societies would not be different from the RMs (i.e., without any threats). To explore this, we generated a heat map illustrating the difference in frequency of occurrence of successful societies within

that specific parameter range (# successful societies within parameter range/ total # of successful societies) and the frequency expected given the number of simulations run within that given parameter space (# simulations run within parameter range/total # of simulations run). Every tile in Fig. 3 represents simulations within a specific parameter space range. These ranges comprise steps of 0.5 and 1 for *lambda threat rate* and *lambda rate intensity* respectively (12–13 in Table 1 and Fig. 3). The color code shows the difference between observed and expected. The yellow-white areas show that the condition was more frequent and orange-red areas that it was much less frequent in successful societies than in the whole parameter sampling. Hence, societies thrive in the white-yellow zone, when the rate of threats is low-medium (i.e., 0.5–7) and the intensity of threats is low (i.e., 30–50, the higher the value of lambda intensity, the lower the intensity of threats) (Fig. 3).

Finally, using the same parameter ranges as in Fig. 3, we generated a heatmap plotting the average religiosity of the successful societies falling within that specific parameter range (Fig. 4). The average religiosity of the society increases the higher the rate (high values x-axis) and intensity (low values y-axis) of threats (Fig. 4).



**Fig. 3** Heat map of the differences between the percentage of successful societies observed and the percentage expected given the number of simulations run within that specific parameter range for RM **a** 1 and **b** 2. Color scale are the difference between % observed–% expected



**Fig. 4** Heat map of average religiosity of successful societies falling within a specific parameter range for **a** RM 1 and **b** RM 2. Color scale are values of average religiosity

## 4 Discussion

The simulation results presented above have shown the plausibility of the central idea underlying the model—that large scale societies may have been made possible by religiously-motivated pro-social behavior. In the model, societies in which the agents readily performed pro-social behavior were able not just to survive when facing environmental threats but even to grow many times beyond their initial size. The relative harshness of the environments faced by the societies can be understood when we consider that in the initial sensitivity analysis only a very small minority of societies was successful.

The conditions that were identified as overwhelmingly present among successful societies across the two reference models give us a lot of additional insight. In particular, it is clear that it is important that agents be readily willing to engage in pro-social behavior, that the behavior not be particularly costly, that it benefits a larger number of individuals, and that insecurity plays a smaller role in whether people have children than their wellbeing. Most significantly, it was shown that: (1) pro-social behavior had to be efficient, i.e., its cost had to be smaller than its effect on security and religiosity; and (2) the less effective the pro-social behavior, the more readily the agents must be willing to engage in it. While these results are not fundamentally surprising, they do show that the modelled societies are behaving in ways that appear to capture important aspects of reality.

The key results for the plausibility of the model, however, were those showing the relationship between threat levels, on the one hand, and the religiosity and success of societies, on the other. Firstly, it was clear that the most religious societies were those facing the most severe and most frequent threats—as has been seen in many historical real-life cases. Interestingly, the rate of threats appears to be more significant—showing that infrequent but large threats are not enough to maintain very high religiosity in a society.

Secondly, the relationship between threat levels and success was more involved. When it came to success, the pattern is similar in that the greatest number of successful societies is to be met where neither the rate nor the severity of threats is too great. However, an interesting difference is that very low threat rates do not lead to the highest rates of success but, instead, seem to be connected with somewhat decreased success. This is most probably because in the intermediate conditions, religiously-motivated pro-social behaviors could counteract the insecurity while also maintaining high levels of religiosity. In environments where the threats were more intense or more frequent, even high levels of religiosity and resulting pro-social behavior were not enough to keep insecurity low and allow the societies to succeed. Unlike religiosity, success appears to be more connected to the average size of the threats—showing that less infrequent but large threats can overwhelm a society that has not maintained sufficiently high levels of cooperation. This is also likely to be connected to the fact that societies facing the lowest frequency of threats were not particularly successful even if those threats were not particularly large.

In future work, we plan to extend the architecture of this Prosocial Equilibrium model in order to address other research questions such as: What is the role of non-religious central institutions (of the sort common in secular societies) in promoting prosociality, lowering anxiety and enhancing wellbeing? These further developments will contribute to major debates in the scientific study of religion and secularization.

# References

1. Hamilton, W.D.: The genetical evolution of social behaviour. II. J. Theor. Biol. **7**(1), 17–52 (1964)
2. Trivers, R.L.: The evolution of reciprocal altruism. Q. Rev. Biol. **46**(1), 35–57 (1971)
3. Williams, G.C.: Adaptation and natural selection. In: Adaptation and Natural Selection. Princeton University Press (2018)
4. Olson, M.: The Logic of Collective Action, vol. 124. Harvard University Press (2009)
5. Fehr, E., Gächter, S.: Altruistic punishment in humans. Nature **415**(6868), 137–140 (2002)
6. Heckathorn, D.D.: Collective action and the second-order free-rider problem. Ration. Soc. **1**(1), 78–100 (1989)
7. Durkheim, E.: The Elementary Forms of Religious Life. Transl. by KE Fields. Free Press, New York (1912)
8. Norenzayan, A.: Big Gods: How Religion Transformed Cooperation and Conflict, 1st edn. Princeton University Press, Princeton (2013)
9. Bulbulia, J.: Religious costs as adaptations that signal altruistic intention. Evol. Cogn. **10**(1), 19–38 (2004)
10. Bering, J., Johnson, D.: 'O Lord… You Perceive my Thoughts from Afar': recursiveness and the evolution of supernatural agency. J. Cogn. Cult. **5**(1–2), 118–142 (2005)
11. Shults, F.L., Wildman, W.J., Lane, J.E., Lynch, C., Diallo, S.: Multiple axialities: a computational model of the axial age. J. Cogn. Cult. **18**(4), 537–564 (2018)
12. Malinowski, B.: Magic. Sci. Relig. N. Y. Doubleday (1948)
13. Norris, P., Inglehart, R.: Sacred and Secular: Religion and Politics Worldwide, 2nd edn. Cambridge University Press, Cambridge (2011)
14. Lang, M., et al.: Moralizing gods, impartiality and religious parochialism across 15 societies. Proc. R. Soc. B **286**(1898), 20190202 (2019)
15. Talmont-Kaminski, K.: Religion as Magical Ideology: How the Supernatural Reflects Rationality. Routledge (2014)
16. Tsang, J.-A., Al-Kire, R.L., Ratchford, J.L.: Prosociality and religion. Curr. Opin. Psychol. **40**, 67–72 (2021). https://doi.org/10.1016/j.copsyc.2020.08.025
17. Gore, R.J., Lynch, C.J., Kavak, H.: Applying statistical debugging for enhanced trace validation of agent-based models. SIMULATION **93**(4), 273–284 (2017)
18. Diallo, S.Y., Gore, R., Lynch, C.J., Padilla, J.J.: Formal methods, statistical debugging and exploratory analysis in support of system development: towards a verification and validation calculator tool. Int. J. Model. Simul. Sci. Comput. **7**(01), 1641001 (2016)
19. Gore, R., Reynolds, P.F., Jr., Kamensky, D., Diallo, S., Padilla, J.: Statistical debugging for simulations. ACM Trans. Model. Comput. Simul. TOMACS **25**(3), 1–26 (2015)

# An Agent-Based Model of the Role of Epistemic Vigilance in Human Cooperation

**Mariusz Rybnik, Ivan Puga-Gonzalez, F. LeRon Shults, Ewa Dabrowska-Prokopowska, and Konrad Talmont-Kaminski**

**Abstract** This paper uses an agent-based model with an adapted stag hunt style scenario to explore the role of the social transmission of correct information about stag hunting and potentially incorrect information about the costs of defection on cooperation in a small artificial society. The computational architecture of the model draws upon Daniel Sperber and Hugo Mercier's concept of epistemic vigilance as well as Brian Skyrms' work on cooperation in stag-hunt scenarios. In the model, communities of 100 hunters begin with no knowledge of stag hunting or punishment for defection and via imperfect social learning, guided by source or content vigilance, move toward a stag hunting or hare hunting equilibrium, where stag hunting may be motivated by the expectation of cooperation or by the fear of punishment. Most successful communities end up using content vigilance to determine their beliefs regarding stag hunting but use source vigilance to determine their beliefs regarding punishment, as predicted in the theoretical work of Konrad Talmont-Kaminski. These findings contribute to the ongoing debate in a variety of disciplines about the conditions under which—and the mechanisms by which—cooperation emerges and is maintained in human societies.

**Keywords** Agent-based model · Epistemic vigilance · Prosocial behavior · Human cooperation

M. Rybnik (✉)
Institute of Computer Science, University of Bialystok, Bialystok, Poland
e-mail: m.rybnik@uwb.edu.pl

I. Puga-Gonzalez · F. L. Shults
NORCE Center for Modeling Social Systems, Kristiansand, Norway

E. Dabrowska-Prokopowska
Institute of Sociology, University of Bialystok, Bialystok, Poland

K. Talmont-Kaminski
Society and Cognition Unit, University of Bialystok, Bialystok, Poland

# 1 Introduction

Here we outline and present the initial experimental results of a computational model of the role of epistemic vigilance in human cooperation using an adapted stag hunt style scenario. The agent architectures for the artificial society are informed by the theory of epistemic vigilance originally formulated by Sperber and Mercier [3], and further articulated by Talmont-Kaminski [4], as explained in more detail below. The stag hunt game was originally proposed in [2] as a formal description of a scenario where the highest utility equilibrium requires mutual trust between the players. Inspired by a Rousseau story [1], it places two hunters in the situation of needing to independently determine whether to hunt a stag together or a hare individually. Hunting the stag offers a significantly higher payoff but requires the cooperation of the other hunter. Hunting the hare offers a guaranteed but lower payoff. The game is similar to the classic prisoner's dilemma, which becomes a stag hunt scenario if defection is punished sufficiently to lower the payoff for defection below that of reciprocated cooperation. As Skyrms has shown, providing hunters can choose whom to hunt with, the norm of cooperation typically becomes stabilised in a population of hunters [2].

In Skyrms' scenario the only unknown is whether the other player will cooperate—all the costs of the alternative strategies are known and mutually cooperating stag hunters are invariably successful. To explore the significance of knowledge regarding those variables it is necessary to draw upon epistemic vigilance theory as set out in [3]. As pointed out there, social learning makes it possible for people to benefit from the experience of others without the potential costs of obtaining that experience but at the cost of the potential for accidental or deliberate deception. To avoid this, people engage in epistemic vigilance by using cues regarding the content of the information being presented to them, as well as the source of that information to determine whether to accept it. In effect, we can consider the plausibility of what someone has told us as well as the trustworthiness of that individual in order to decide whether to believe what was said.

The choice between using content versus source vigilance has important social consequences. Talmont-Kaminski argues in [4] that, where the utility of a belief is tied to its accuracy, a preference for content vigilance will lead to better results at the social level whereas a preference for source vigilance is to be favoured where utility is not connected to accuracy. The same does not necessarily hold at the individual level. As a prime example of beliefs whose utility is not dependent upon their accuracy he proposes beliefs that motivate cooperation. This opens the way to combining the stag hunt with epistemic vigilance to empirically examine how the choice between source versus content vigilance regarding information about potential punishment and stag behaviour interacts with a repeated many-player stag hunt scenario that allows for partner choice. The results are informative regarding the potential role of the fear of supernatural punishment for motivating cooperation.

# 2 ABM Stag Hunting Variant

The model represents a small-scale society of a hundred hunters who initially do not cooperate, have no knowledge of how to hunt stags and have no belief in punishment regarding defection. The hunters set out in groups of five. Group of five individuals is the minimal group size that allows to introduce social cooperation to the original stag game, with eventual great reward if majority of group cooperated. Each hunter decides to cooperate either because of previous cooperation by the group members or because due to punishment they believe the utility of defection will be lower than that of an unsuccessful hunt. A minimum of three hunters must cooperate for the hunt to have any chance to be successful, its odds of success depending upon the highest stag behaviour knowledge of them. Between hunts, hunters exchange two pieces of information, one about stag behavior and the other about belief in punishment. For each piece of information, they use source or content vigilance to determine whether to accept it or not. Additionally, every ten hunt attempts, hunters evaluate their relative hunting success and whether to modify their strategy of relying upon source or content vigilance. There is no real punishment and no individual learning regarding either stag behavior or punishment beliefs.

## 2.1 Agent Variables

- *Competence CMP; CMP*$\in$ [0, 1]—agent's ability to perform content vigilance, i.e. judge the accuracy of beliefs. Ranges from 0 to 1—no competence to perfect competence. Generated randomly at start using uniform distribution. Does not change.
- *Status STS; STS*$\in$ [0, 1]—agent's value in considerations of source vigilance, i.e. how they and others estimate their value as source. Ranges from 0 to 1—the lowest status to the highest status. Generated randomly at start using uniform distribution. Does not change.
- *Stag beliefs SB; SB*$\in$ [0, 1]—agent's beliefs regarding stag behaviour and necessary to catch the stag. Ranges from 0 to 1—completely misleading knowledge to perfect knowledge. Starts at 0, but changes due to interactions with other agents.
- *Punishment beliefs PB; PB*$\in$ [0, 1] >—agent's beliefs regarding the size of punishment for defection. Ranges from 0 to 1—completely misleading knowledge to perfect knowledge. Initially 1 for all agents, representing the perfect knowledge that there is no punishment in the model. Changes due to interactions with other agents.
- *Cooperated/Defected CD; CD*$\in$ {0, 1}—agent's behaviour in the previous hunt. 0 if defected, 1 if cooperated. Set after first hunt on basis of initial choice. Updated after every hunt.

- *Stag belief vigilance preference—SVP; SVP*∈ {*source, content*}—agent's preference for source or content vigilance when evaluating stag beliefs. 1 = content preferred, 0 = source preferred.
- *Punishment belief vigilance preference—PVP; PVP*∈ {*source, content*}—agent's preference for source or content vigilance when evaluating punishment beliefs. 1 = content preferred, 0 = source preferred.
- *Utility Total UT*—a variable representing the total utility obtained by the agent during current set of 10 hunts. Set to 0 at the start of every set of hunts and updated after every hunt. Strategies are updated each 10 hunts using the *UT* in order to identify more successful individuals and allow the less successful ones to adopt strategies from them.

## 2.2 Global Variables

- *Group assortment GA—GA*∈ [0; 1]—ability of cooperating agents to pick groups with cooperating agents if given the opportunity. Ranges from 0 to 1—none to complete. Determines the degree to which co-operators will group together.
- *Village size*—number of agents in the village. It remains constant at 100 agents. No deaths/births occur in the model.
- *Cooperate/defect error CDE—CDE*∈ [0; 1]—probability of agents incorrectly remembering behaviour of other agents in previous hunt. Ranges from 0 to 1—no error, to always error.
- *Strategy update group SUG—SUG*∈ [1; 50]—the percentage of agents who update their strategy regarding use of source or content vigilance. It ranges between 1 and 50 percent of the village size (n=100).
- *Social Learning Ratio SLR—SLR*∈ [0.1; 1]—ratio of learning beliefs (*SB* or *PB*) during *Social Learning* phase. When agents consider worthy to learn the beliefs of another agents, their beliefs moves towards that of the partner with a ratio given by *SLR*.
- *Social Learning Error—SLE*∈ [0.01; 0.30]—absolute value of error when learning the beliefs of another agent (*SB* or *PB*). The learning occurs during *Social Learning* phase.

## 2.3 ABM Cycle (Steps)

The model's flowchart is presented in Fig. 1, with the sequential steps descriptions following.

**Step 1. Assignment to hunting groups**. At initialization, agents are assigned into groups of five randomly. Afterwards, before each hunt agents are anew assigned into groups in the following order: Agents who defected in previous round—assigned

**Fig. 1** The model flowchart and direct influences (represented with colorful arrows). Please note that Hunter (encircled in blue) represents a single agent. *Competence* and *Status* are crossed as being constant hunter properties. 5 external model parameters are to the right and denoted with circle-with-arrow symbols with arrows showing the step they affect

to a random group with free spots. Agents who cooperated in previous round—if a random number between 0 and 1 is below *GA* value, the agent is assigned to a group (when available) with free spots that contains hunters who cooperated in the previous round. Otherwise, the agent is assigned to a random group with free spots.

**Step 2. Choose to cooperate/defect**. In the first hunt, no agent cooperates. In later hunts, all agents determine whether they will cooperate or defect in a given hunt on the basis of:

- Expected utility of defection $XUD = -1 + (PBx2)$. Ranges from -1 to 1.
- Expected utility of cooperation $XUC =$ for each agent in the group (including self) add $CD$ (with $CDE$ being probability of misremembering own/group member $CD$). With $CD$ taking a value of 1 if the agent cooperated in the previous hunt or 0 if not. If however dis-remembered ($CDE$ fired) respectively taking the opposite value. Hence, $XUC$ ranges from 0 to 5—none to all agents are thought to have cooperated in the previous hunt.
- Level of *SB* is ignored.

**Step 3. Determine outcome**. All defectors increase $UT$ by 1. They caught their hare and were not punished for it as there is no punishment in the model.

All groups with less than three cooperators do not change $UT$. There were not enough cooperators to catch the stag.

All groups with three or more cooperators increase $UT$ by 5 with probability equal to the highest $SB$ among them. The group was large enough to potentially catch the stag on the basis of the best knowledge available to them.

```
if XUD < 0 then
  │ cooperate ← true ;        /* It is better to catch nothing than to
  │ catch a rabbit and be sorely punished. */
else
  │ if XUC < 3 then
  │   │ cooperate ← false;    /* Assuming this group behaves the same
  │   │ as last time, there will not be enough hunters to catch
  │   │ the stag. */
  │ else
  │   │ cooperate ← true;
  │ end
end
```

**Algorithm 1:** Choose to cooperate/defect

**Step 4. Social learning**. After each hunt, the villagers exchange and potentially update their beliefs. They determine whether to update their beliefs on the basis of either source or content vigilance (depending upon their preferences). In the first case, they compare their status with that of the villager they are communicating with. In the second case, they use their competence to estimate the accuracy of the beliefs in question and learn only when it seems profitable.

Each villager:

- partners up with another random villager,
- can partially adopt $SB$ and $PB$ from their partner, if the condition based on corresponding vigilance preference is fulfilled (as specified in details below). The degree of adopting beliefs is defined as the parameter *Social Learning Ratio (SLR)*.

On the basis of their $SVP$ and $PVP$ the villager determines whether to update their beliefs:

- If preferring source vigilance, use status *(STS)*:
  - Calculate status difference $STSD$ = partner's $STS$—Ego's $STS$
  - *STSD* determines probability of shifting towards partner's belief with the probability ranging from 0 when $STSD$ = -1; 0.5 when $STSD$ = 0; and 1 when $STSD$ = 1.
  - If the condition is fulfilled, own belief will be shifting towards the partner with a small ± error (the parameter *Social Learning Error*).
- If preferring content vigilance, use competence *(CMP)*:
  - Calculate accuracy for partner's belief and ego's belief using *CMP*. The higher the competence, the more precise the assessment—random value if $CMP$ = 0, the precise actual value if $CMP$ = 1.
  - Shift towards partner's belief with ± error (the parameter *Learning Ratio Error*), if accuracy of partner's belief is estimated to be higher than of own belief.

**Loop** Steps one through four (from forming hunting groups to determining who cooperates, outcomes of the hunts, and social learning) are repeated in a loop ten times, allowing the agents to collect utility.

**Step 5. Update strategy**. After ten rounds of hunting and communicating, the agents evaluate the strategy they have been using to update their beliefs: source or content vigilance. Agents are ordered by their UT and each of those in the percentage (*SUG*) of agents with the lowest UT:

- Randomly choose between *SVP* and *PVP*
- Randomly choose an agent from those with very high UT (use *SUG* to determine size of group).
- Adopt the vigilance preference (i.e. *SVP* or *PVP*) of that agent.

Finally, all agents set their *UT* back to 0 and the cycle starts again.

**Desired outcome of the game variant** It is expected that over time agents will come to select stag beliefs on the basis of content vigilance and punishment beliefs on the basis of source vigilance resulting in societies dominated by highly accurate beliefs about stags and highly inaccurate beliefs about punishment. The model is specifically designed so that the two kinds of beliefs are treated exactly the same way but differences may arise because of the way the beliefs impact hunting success. This is so that any differences in source/content vigilance regarding them can be traced back to this practical difference.

## 2.4 The Model

The model was implemented in AnyLogic v.8.7.10. Here we present a brief description of the model. A full ODD+D protocol description can be found at the repository: https://github.com/mrybnik/staghunters.

## 3 Results

### 3.1 General Model Behaviour

**Initial state** The agent population starts with perfect knowledge of the lack of punishment in the model (*PB* = 1). The population start with no knowledge of stag behaviour (*SB*=0). The *SVP* and *PVP* strategies are varied in population, both statistically close to half *content* and *status*.

**Hare hunting equilibrium** The population starts with no cooperation and zeroed stag beliefs. While stag beliefs gradually increase due to learning error and then spread in the population due to social learning, cooperation initially remains low.

**Fig. 2** Emerge of spontaneous cooperation. Blue area represents averaged population cooperation. Both brown and red lines are quite high, representing averaged *SB* and *PB*. Brown area presents *SVP* preferences, being mostly *source*-based. Red area presents *PVP* preferences, being mostly *content*-based. Finally green area represents averaged *Total Utility* with a quick raise from *hare hunting equilibrium* to mostly successful *stag hunting equilibrium*

Almost all hunters hunt hares and stag hunting is marginally popular. This could be called *Hare hunting equilibrium*.

**Fear-induced cooperation** False beliefs regarding punishment may appear due to learning error and spread due to social learning, especially with *Punishment Beliefs* learning driven by source *Vigilance Preference* (i.e. *PVP=source*). Low punishment beliefs lead to cooperation due to the fear of being punished when not cooperating, even though in reality punishment does not exist. Cooperation—once ignited with fear—spreads in population, resulting in almost all hunters cooperating in hunting stags—so called *Stag hunting equilibrium*.

**Spontaneous cooperation** Mistaken recall of prior cooperation (*CDE*) may lead to spontaneous cooperation, which may sometimes spread within the population resulting in the *Stag hunting equilibrium*. An example of this phenomenon is depicted in Fig. 2. The effect is enhanced by high values of *Group Assortment* parameter (*GA*).

**Stag hunting equilibrium** Once cooperation appears, be it due to fear or spontaneous propagation, the whole population generally cooperates all the time. This can be called *Stag hunting equilibrium*. Providing the *SB* is high enough, this results in *Utility*

*Total* being close to 50. In ideal conditions (perfect recall of previous cooperation *CDE*=0 and very high *Stag Beliefs*) this would result in every group hunting stags successfully.

## 3.2 Sensitivity Analysis: The Role of GA and CDE for Spontaneous Cooperation

A sensitivity analysis was performed using LHS hypercube sampling resulting in 10000 parameter combinations. Sampled parameters:

- *Strategy Update Group SUG* [0.01–0.25],
- *Social Learning Error SLE* [0.01–0.2],
- *Cooperate Defect Error CDE* [0.01–0.3],
- *Group Assortment GA* [0.01–1],
- *Social Learning Ratio SLE* [0.01–1].

Figure 3 presents the result of the analysis with percentage of cooperating agents (y-axis), level of group assortment (x-axis), different ranges of cooperation/defection error *CDE* (facets), and Vigilance Preferences (SVP-PVP) followed finally by the whole population (color and shape of data points).



**Fig. 3** Sensitivity analysis, 4 final combinations of *SVP* and *PVP* are shown

The conclusions here are:

- Fear induced cooperation (with very high cooperation values) is visible for facets with *CDE* in range [0–0.15] and almost exclusively connected with source vigilance regarding punishment beliefs (Fig. 3: *Vigilance Preferences* **content-source** and **source-source**, blue squares and red dots)—which is necessary for false beliefs about punishment to spread.
- Villages that do not fear punishment (Fig. 3: *Vigilance Preferences* **content-source** and **source-source**, blue squares and red dots) with low values of *CDE* [0.0–0.1], regardless of *Group Assortment* value very rarely induce cooperation spontaneously. Influence of *Group Assortment* is however most visible with *CDE* in range [0.1–0.2] where it most probably helps to induce spontaneous cooperation. One can observe that raising *GA* values clearly incite cooperation in this case.
- With *CDE* over 0.3 cooperation becomes erratic for villages that apply content vigilance to both stag and punishment beliefs (Fig. 3: *Vigilance Preferences* **content-content** and **source-content**, violet crosses and green triangles), as agents cannot remember correctly who cooperated recently. This is not a problem for villages using source vigilance regarding punishment beliefs as they rely on fear-induced cooperation rather than cooperation based upon trust (Fig. 3: *Vigilance Preferences* **content-source** and **source-source**, blue squares and red dots).

### 3.3   The Influence of Stag Beliefs and Cooperation on Utility

Figure 4 presents the distribution of UT vs the proportion of CD in the villages for the four combinations of *Vigilance Preferences* at the end of the simulation.

Conclusions: High final *Utility Total* requires both cooperation and high stag beliefs (proficiency in stag hunting). The obvious positive correlation can be observed when content vigilance is applied to stag beliefs (facets *(SVP=content, PVP=source)* and *(SVP=content, PVP=content)* in Fig. 4), where stag hunting proficiency is learnt efficiently and therefore high. Where source vigilance regarding stag beliefs is used (facets *(SVP=source, PVP=source)* and *(SVP=source, PVP=content)*, Fig. 4), however, high cooperation actually tends to worsen the population's *Utility Total*, as attempts to hunt stags are unsuccessful, due to low *SB* (hunting proficiency).

## 4   Discussion

The model created is complex as agents constantly change most of their properties and therefore behavior. The agents evolve through social interactions. Thus it may be seen as evolutionary model rather than a typical ABM, where agents usually behave in a constant manner. While full identification and prediction of all processes seems to be impossible, it is possible to identify and statistically confirm core authors'

**Fig. 4** UT vs CD, 4 combinations of *Vigilance Preferences* at the end of the simulation

assumptions. While of course simpler model could be used, the ultimate goal is to simulate social interactions and try to understand them.

The aim of the model was to examine the role of the social transmission of potentially incorrect beliefs upon cooperation and hunting success in a stag-hunt like scenario. In the case of knowledge of stag hunting it was found to be preferable for a society to rely upon content vigilance for determining whether to accept socially transmitted information. The situation, however, was different in the case of beliefs regarding punishment for failing to cooperate with other hunters. Here, in some scenarios, it turned out to be preferable for members of a society to incorrectly believe in such punishment as it could both induce and maintain very high levels of cooperation—the necessary levels of false belief in punishment requiring a focus upon source vigilance in the case of punishment beliefs. In particular, where the chance of mis-remembering prior behaviour was low—making spontaneous cooperation unlikely—source vigilance regarding punishment beliefs helped to induce

cooperation while, where prior behaviour was often mis-remembered—leading to inconsistent trust in other hunters—source vigilance regarding punishment beliefs helped to maintain much higher steady levels of cooperation. Only for a small range of values for recall of prior behaviour and ability to use that to choose partners did levels of cooperation maintained by trust equal those maintained by the fear of nonexistent punishment.

In his study of human cooperation, Brian Skyrms assumed that cooperation was based upon accurate knowledge of the situation and rational decisions regarding expected utility. We have expanded upon his work by considering the potential effect of incorrect information passed on by social learning, driven by epistemic vigilance towards the source or content of that information. In effect, we have shown that a preference for source vigilance—while rational at the individual level—has crucial significance for the spread of incorrect information with completely different consequences at the social level depending upon whether the utility of that information being connected to its accuracy. And in the case of information that motivates cooperation, utility and accuracy can be quite separate—as proposed by Talmont-Kaminski.

## References

1. Rousseau, J.J., et al.: Discours sur l'origine et les fondements de l'inégalité parmi les hommes, suivi de la reine fantasque. Discours sur l'origine et les fondements de l'inégalité parmi les hommes, suivi de La reine fantasque, pp. 1–190 (2009)
2. Skyrms, B.: The Stag Hunt and the Evolution of Social Structure. Cambridge University Press (2004)
3. Sperber, D., Clément, F., Heintz, C., Mascaro, O., Mercier, H., Origgi, G., Wilson, D.: Epistemic vigilance. Mind Lang. **25**(4), 359–393 (2010). https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1468-0017.2010.01394.x
4. Talmont-Kaminski, K.: Epistemic vigilance and the science/religion distinction. J. Cogn. Culture **20**(1–2), 88–99 (2020). https://brill.com/view/journals/jocc/20/1-2/article-p88_5.xml

# Attracted to Fish: A Gravity-Based Model of Purse-Seine Vessel Behaviour

**Nicolas Payette, Ernesto Carella, Katyana Vert-Pre, Brian Powers, Steven Saul, Michael Drexler, Aarthi Ananthanarayanan, and Richard Bailey**

**Abstract** This paper presents a gravity-based behavioral algorithm designed to simulate the dynamic decision-making processes of purse seine fishers in the Eastern Pacific Ocean. The algorithm captures the complex interplay between fishers' actions, environmental conditions, and regulatory constraints. It comprises two core strategies: an action strategy and a destination strategy. The action strategy involves selecting the most favorable course of action based on estimated values and preferences, while the destination strategy uses gravity fields to determine attractive ocean cell locations. These fields are modulated by real-time circumstances, guiding fishers toward areas of high value. Calibration against real-world data from the Inter-American Tropical Tuna Commission (IATTC) observer database is ongoing, with a focus on achieving accurate representation of action frequencies and species-specific catch per action type. Initial calibration results highlight the need for further refinement. While still a work in progress, this algorithm provides a robust foundation for capturing the intricate dynamics of purse seine fishing, adapting to evolving conditions, and informing policy evaluations. Future enhancements include adaptive fishing strategies and incorporating fleet-level interactions for a more comprehensive understanding of fishing behaviors.

**Keywords** Fishing behavior · Fish aggregation devices · Gravity model · Agent-based simulation · Fisheries management

N. Payette (✉) · E. Carella · R. Bailey
University of Oxford, Oxford, UK
e-mail: nicolaspayette@gmail.com

K. Vert-Pre · B. Powers · S. Saul
Arizona State University, Tempe, USA

M. Drexler · A. Ananthanarayanan
Ocean Conservancy, Washington, D.C., USA

Fisheries and agent-based models are a match made in heaven [2]. Fisheries are complex human-environmental systems involving adaptive agents that interact with each other and with their surroundings in a spatially explicit environment—a textbook definition of what ABMs are good at. Fish and other seafood represent about 17 percent of global animal protein consumption [7] so ensuring ecological and economic sustainability of fisheries is vitally important but it is difficult to predict the impact of policy interventions on such complex systems.

The POSEIDON agent-based model [1] aims to ease these difficulties for scientists and policy makers trying to compare different possible scenarios. The model, built in Java using the MASON library [12], is a flexible tool kit that can be applied to a wide range of fisheries. Applications have so far included US West Coast Groundfish [4], Indonesian deepwater snapper/grouper mixed fishery (forthcoming) and more conceptual work [3, 13].



**Fig. 1** Model layers and interactions

Our latest ongoing project tackles the tuna purse-seine fishery in the Eastern Pacific Ocean (EPO), an area managed by the Inter-American Tropical Tuna Commission. This fishery is characterised by intensive use of fish aggregating devices (FADs). It involves three main species of tuna: Bigeye (*Thunnus Obesus*), Yellowfin (*Thunnus Albacares*) and Skipjack (*Katsuwonus Pelamis*). Bigeye is vulnerable and regulations aim to protect it.

The POSEIDON EPO model consists of multiple interacting layers (Fig. 1), each of which could be the subject of a whole paper. Here we focus on the fleet behaviour model which is still very much a work in progress. We will give a high-level description of the challenges we face in modelling the EPO purse seine vessel fleet and how we approach them using a "gravity-based" approach where agents make their decisions based on the perceived value of different locations in the environment. Given the state of the work, we offer no strong conclusions at this point. Still, we believe our approach to be novel and interesting for the social simulation community to learn about.

# 1   The Eastern Pacific Tuna Fishery

Tunas like to congregate under floating objects like drifting logs or algea. The reasons for this behaviour are not well understood [14]. It could be that floating objects act as social meeting points [6] or that they are generally indicative of a resource rich environment [8].

Whatever the reasons for it, fishers have long been taking advantage of this behaviour by using "Fish aggregating devices" (FADs): human-made raft-like assemblages that can be monitored for the presence of fish and targeted when the aggregation gets big enough. Fishers use purse seines to catch tuna: large nets, with floats at the top and weights at the bottom, that can surround an entire school of fish. Catching fish using this method is called "making a set".

FADs in the EPO are often left to drift over long distances and can attract hundreds of tons of tuna along the way. The advent of GPS and echo-sounder technology has drastically increased the efficiency of FAD fishing and it is now the main tuna fishing method in the EPO. It is not the only one, however. Sets are sometimes made on "dolphin-associated" schools of tuna, as mature tunas (especially Yellowfin) tend to band together with dolphins. Dolphins are sometimes killed in the process, and only a limited number of vessels are permitted to make dolphin sets. Finally, some sets are "non-associated": they target free schools of tuna with neither dolphins nor floating objects involved.

While most fishers working with FADs rely on deploying their own FADs and then tracking their position and productivity using satellite buoys, it can happen that a vessel encounters a FAD that is not theirs. If tuna is present at the FAD, it is there to be taken. The ethics of making these opportunistic FAD sets are controversial amongst fishers, but the practice is legal and common enough that it warrants inclusion in the model.

The IATTC has an on-board observer program and collects data about the different sets performed by vessels, including their location and how much fish they catch. By using the clustering analysis method described in [9], they can identify the fishing strategies practised by different parts of the fleet. Recent applications of this analysis have identified three main strategies [11]. Cluster A, roughly a quarter of the fleet, relies on dolphin sets. They make occasional FAD sets and non-associated sets, but deploy few FADs of their own. Cluster B, also about a quarter of the fleet, relies almost exclusively on FADs that they have deployed themselves. They tend to be the larger vessels and are prepared to travel further west in pursuit of their FADs. Cluster C, the last half of the fleet, stays closer to the coast and relies on a mixture of non-associated sets, sets on their own FADs and opportunistic FAD sets. Those strategies are fluid and fishers adapt to circumstances.

## 2 The POSEIDON EPO Model

The model represents 181 purse seine fishing vessels of IATTC size class 6 (i.e., the biggest vessels, with over 363 tonnes of carrying capacity each), which account for the bulk of the catch in the EPO. Smaller purse-seiners (class 1–5) and long-line fishing vessels are exogenous to the model: we remove their annual catch from the ocean without simulating them as agents.



**Fig. 2** The EPO as viewed in POSEIDON. The grey areas represent the El Corralito spatial closure area and the Galapagos Islands protected area. The black anchors are ports

The simulated spatial extent (Fig. 2) spans from 50° S, 171° W to 50° N, 70° W, and space is partitioned within this extent into 1° × 1° square grid cells approximately 111 km high and varying between approximately 72 km wide (at latitudes ± 50°) and 111 km wide (at the equator). Within this general area, the area under IATTC regulations (known as the "Antigua Convention Area") extends from longitude 150°W to the west coast of the Americas.

Current regulations in the EPO include a marine protected area (MPA) around the Galapagos Islands, where no fishing is allowed at any time, a yearly spatial closure (September 29th to October 29th) known as "El Corralito", a seasonal closure from either July 29th to October 8th or November 9th to January 19th (at the vessel's discretion), and limits on the number of active FADs that a vessel can have at the same time.

Simulated vessels in POSEIDON exist in discrete space: they move from grid cell to grid cell, but their position is no more precise than that. POSEIDON's biological layer uses the same discretization: each grid cell contains a vector of population per age group for each fish species in the model. Simulated FADs, however, exist in continuous space: each FAD has a precise longitude/latitude position. We use oceanic current vectors derived from the HYCOM model [5] to model the drifting of FADs.

The spatial distribution of fish in the ocean is based on monthly habitability maps produced using the methods described in [10]. As FADs drift, some fish is transferred from the ocean cells to the FADs. The probability of a given FAD to attract fish at any point in time depends on both the quantity of fish in its current ocean cell and the quantity of fish currently over the FADs.

The model assumes that fishers have real time information about the location and quantity of fish under the FADs that they have deployed. This reflect their real-world use of GPS and echo-sounder buoys, though readings from the latter are only approximate in reality.

There are two parts to the fishers' strategy: the *action* strategy (What do I do here?) and the *destination* strategy (Where do I go next?). We will tackle each in turn.

## 3   The Action Strategy

Once a fisher arrives in a cell, they can perform various actions (see Table 1). Which actions are available depends on the current circumstances. In order to deploy a FAD, the agent must have at least one FAD left in their inventory. In order to make a set, an opportunity must exist. In the case of making a set on their own FAD, this simply requires the FAD to be present in the same ocean cell, as the fisher knows the precise locations of their FADs. In order to make an opportunistic set (OFS) on someone else's FAD, the FAD must not only be there, it must also be detected by the fisher. The same goes for dolphin sets (DEL) and non-associated sets (NOA): a school of fish must be present *and* detected. The probabilities of detecting various

**Table 1** Possible action types

| Code | Action |
| --- | --- |
| DPL | Deploy a FAD |
| FAD | Make set on own FAD |
| OFS | Make opportunistic FAD set |
| NOA | Make non-associated set |
| DEL | Make dolphin-associated set |
| SSO | Search for set opportunities |

set opportunities are parameters of the model. Agents also have the option to "Search for set opportunities" (SSO), which increases those probabilities.

Sets on own FADs and opportunistic FAD sets can be represented more than once in the list of available actions if, respectively, the fisher has multiple FADs in the cell or multiple other FADs have been detected. For NOA and DEL sets, we assume that only one school of each type is available to set on at a time (though new opportunities may arise right after a set is made).

Before choosing an action, the fisher must identify the subset of actions that are not only possible, but also *allowed*. Whether or not an action is allowed depends on the active regulations at the time. Seasonal closures, protected areas, and limits on active FADs or number of sets can forbid actions like FAD, OFS, NOA, DEL and DPL. Dolphin sets also require the vessel to be licensed to fish on dolphins.

Once they know which actions are possible and permitted, the fisher must compute a value for each action. The action with the highest value will be executed first, as long as it's above a given "moving threshold" (another parameter of the model). There are three steps to computing an action's value:

1. Estimate the base value of the action.
2. Apply normalization functions to scale values to [0, 1] interval.
3. In the case of deployments and search actions, apply decay.
4. Weight the values according to fisher preferences.

Let's designate the computed value of an action after the first step as $v_1$. The base value of a potential set is simply the monetary value of the fish that's aggregated under the FAD or part of the tuna school, and we assume that fishers can accurately estimate this value. For deployment and search actions, the base value depends on the location. We'll come back to how those values are established when we describe the destination strategy, but the general idea is that the value of a search action depends on the NOA and DEL sets that were historically made in this location by the fisher, and the value of deploying a FAD depends on sets that were historically made on FADs deployed in this location.

Once we have those basic action values, we need to map them all to the [0, 1] interval so that they are comparable to each other. We do this using a logistic function of the form

$$v_2 \leftarrow \frac{1}{1 + e^{-k_a(v_1 - m_a)}}, \tag{1}$$

where $v$ is the value of the action and $k_a$ and $m_a$ are parameters of the model respectively giving us a steepness and a midpoint for each type of action $a \in$ {DPL, FAD, OFS, NOA, DEL, SSO}. Varying those parameters for different types of actions tells the fisher how valuable an action needs to be before it's worth doing. Furthermore, using logistic functions instead of simple threshold functions allows us to keep the ordering of actions that are otherwise close in value.

Next, we have two special cases: deployments and search actions. The value of doing those actions depends on how many times they have just been done in the current location. It might be worth searching for set opportunities a few times, but you have to give up eventually. As for deployments, it is no use deploying all your FADs in the same ocean cell. We handle this by applying exponential decay to the value of the action:

$$v_3 \leftarrow \begin{cases} v_2 \cdot e^{-\lambda_a n_a}, & \text{if } a \in \{\text{DPL, SSO}\} \\ v_2, & \text{otherwise,} \end{cases} \tag{2}$$

where $n_a$ is the number of times the agent has already taken the action and $\lambda_a$ is an action-specific decay rate.

At this point, every action has a value between zero and one. The final step is to weight them by the fisher's preference for different types of actions so that, e.g., dolphin-setting vessels from cluster A value DEL sets more highly. To weight the preferences of the fishers for each type of action, we simply look at the number of events of each type in the IATTC observer database and assign weights to the fishers in proportion of those (see Fig. 3). This captures the difference in behaviour between dolphin-setters and FAD setters, and between the fishers trying to fish on their own FADs versus those adopting a more opportunistic strategy.

The weights then act as coefficients, so the final action value is given by

$$v_4 \leftarrow v_3 \cdot w_{i,a}, \tag{3}$$



**Fig. 3** Action preference weights for modelled vessels

where $w_{i,a}$ is the weight given by individual fisher $i$ to an action of type $a$.

The agent executes all possible actions whose value is greater the moving threshold in descending order of value. Once there are no more such actions, it is time for the fisher to pick a new destination.

## 4 The Destination Strategy

This is where the concept of a gravity-based model really comes into play. The general idea is that each agent maintains a set of "gravity fields" keeping track of how attractive different ocean cells are for different types of actions. The values in each of those fields are modulated by circumstances and then used to compute an attraction vector relative to the fishers current location. Those different vectors are then weighted by fisher preferences and combined into a single vector that tells the agent in which direction to move. Figure 4 summarises the process.

Note how we are not picking a final destination: just a direction of travel. Once the agent arrives in the next cell, it performs the actions that are worth performing there and then recomputes its travel direction vector. This process eventually leads fishers to areas of highest values but is also highly responsive to changing circumstances like FADs drifting or the vessel's hold getting full.

We consider six different sources of attraction, corresponding to some of the basic actions that fishers can take (Table 2).

We use the action codes listed above to index a set of gravity fields $G$, such that $G^{\text{FAD}}$ refers to a matrix of FAD values and $G^{\text{FAD}}_{x,y}$ refers to the value of FADs in the cell $x, y$, where $x \in \{0, \ldots, m-1\}$ and $y \in \{0, \ldots, n-1\}$, with $m$ and $n$ respectively denoting the width and height of our ocean grid in cells.

Each of the gravity fields derives its values from different sources, corresponding roughly to the expected value of taking an action in a particular location. Note the presence of a PRT action, which has a special field representing the vessel's incentive for returning to port as the trip gets longer and its hold gets fuller. Despite the fact that most of these field's values are monetary values, they need not all be on the same scale. When choosing between potential destinations, the attraction vector generated

**Table 2** Attraction fields

| Code | Value source | Updated |
|---|---|---|
| FAD | Fish under FADs in cell | As FADs drift, attract and lose fish |
| OFS, NOA, DEL | Previous sets in cell | As sets are made + exponential decay |
| DPL | FADs deployed there that led to sets | As sets are made + exponential decay |
| PRT | 1 at port, 0 elsewhere | Never (for now) |

| | |
|---|---|
| Value fields | Compute cell values for `FAD`, `OFS`, `NOA`, `DEL`, `DPL` and `PRT` fields. |
| Modulate | Adjust the cell values according to circumstances |
| Transform | Compute the gravity vectors from the modulated values |
| Combine | Add the resulting vectors, weighted by fisher preferences |
| Direction vector | Head to the cell in that direction |

**Fig. 4** The destination strategy

from each gravity field will be normalized so that it has a length of one. The resulting vectors will also be weighted according to agent preferences.

But more on that later. For now, let's see how each field gets their values.

## 5 FAD Sets

The $G^{\text{FAD}}$ field is the most dynamic field since its values change as the FADs attract and release fish and drift from cell to cell.

Let $F_{x,y}$ denote the set of FADs currently in ocean cell $x$, $y$ and $v(f)$ be a function that computes the monetary value of the fish currently aggregated under FAD $f$. The current value of a cell is the total value of FADs in that cell:

$$G^{\text{FAD}}_{x,y} = \sum_{f \in F_{x,y}} v(f). \tag{4}$$

Since fishers know the position of their own FADs and their aggregated biomass at all time, the field requires no prior initialization: the values are always computed on the fly.

## 6  Opportunistic FAD Sets, Non-associated Sets and Dolphin sets

When it comes to other types of sets that can be made, fishers do not have real-time information that they can use to establish the value of ocean cells, so they must rely on historical information.

We have observer data about sets that were made by fishers in the real world so we use those to initialize the value matrices of agents in the simulation.

Let $S_{x,y}^a$ be the (mathematical) set of (fishing) sets historically made by a fisher in ocean cell $x$, $y$, where $a$ is the type of fishing set. Let us also suppose another function $v(s)$ giving the monetary value of a fishing set $s$. We initialize the value matrices for $a \in \{$OFS, NOA, DEL$\}$, $x \in \{0, \ldots, m-1\}$ and $y \in \{0, \ldots, n-1\}$ as:

$$G_{x,y}^a \leftarrow \sum_{s \in S_{x,y}^a} v(s). \tag{5}$$

In other words, the value of a cell for sets of a particular type is the total value of the sets of that type that were recorded in that location for the fisher in the observer data (example in Fig. 5).

As the simulation is running and agents are made aware of opportunities for sets of different types, they update the corresponding matrix location $G_{x,y}^a$ by adding the value of the observed potential fishing set $s$:

$$G_{x,y}^a \leftarrow G_{x,y}^a + v(s). \tag{6}$$

We say "potential" because agents take the information into account even if they decide not to act on the opportunity.

To reflect the fact that the environment is changing and that recent information should be weighted higher than previous information, we apply exponential decay to these matrices at the end of each simulated year. Each action type $a \in \{$OFS, NOA, DEL$\}$ has its own decay rate $\Lambda_a \in [0, 1]$ and the matrices are updated by applying $G^a \leftarrow \Lambda_a G^a$. These decay rates are parameters of the model and common to all agents.



**Fig. 5**  Example location values for different set types

## 7  FAD Deployments

The FAD deployment values matrix is similar to the set value matrices, in the sense that it is initialised from empirical data, but those initial values are calculated differently.

While sets have a value directly associated with them (i.e., how much the fish that can be caught from that set is worth), FAD deployments do not immediately reveal their value. That value is known only when the deployed FAD is eventually the object of a set.

The FAD sets recorded in the observer data do not provide the identity of the FAD that the set was made on so, contrary to what's happening in the simulation, it's impossible to retrace the deployment that eventually led to the set. Since we cannot associate specific sets to specific FADs, we take instead the total value of all sets made by a fisher on their own FADs and redistribute that value equally to all deployments, and then add up the values of the deployments in each cell. The matrix is initialized with:

$$G_{x,y}^{\text{DPL}} \leftarrow |D_{x,y}| \sum_{s \in S^{\text{FAD}}} \frac{v(s)}{|D|}, \tag{7}$$

where $|D_{x,y}|$ is the number of deployments made by the fisher in cell $x, y$ and $|D|$ is the total number of deployments made by that fisher, with $S^{\text{FAD}}$ being the set of FAD sets made by that fisher.

When a simulated fisher makes a set on one of their own FADs, the deployment value matrix is updated for the cell $x, y$ where the FAD was originally deployed (*not* the cell where the set happened). Again, we simply update the matrix at that location by adding the value of the set:

$$G_{x,y}^{\text{DPL}} \leftarrow G_{x,y}^{\text{DPL}} + v(s) \tag{8}$$

The $G^{\text{DPL}}$ matrix is also subject to decay and gets updated at the end of the year with $G^{\text{DPL}} \leftarrow \Lambda_{\text{DPL}} G^{\text{DPL}}$.

## 8  Port Attraction

The $G^{\text{PRT}}$ matrix is a single-entry matrix where the location of the home port of the fisher gets the value 1 and all other locations get the value 0.

The port value matrix never needs to be updated but, as we will see shortly, its modulation function adjusts the attraction of the port depending on hold fullness and current trip duration.

The port attraction field is also a special case when it comes to actions weights: for all fishers, $w_{\text{PRT}} = 1$, such that the weight of the port attraction is equal to the

weight of all other fields, ensuring that returning to port is treated as a priority when it's time to do so.

## 9   Field Values Modulation

The modulation functions allow the values in the different fields to be adjusted according to circumstances. Those functions are of the form $\mu(\mathbf{c}, t_{\mathbf{c}}, \Omega)$, where $\mathbf{c}$ is the cell whose attraction we are calculating and $\Omega$ represents the general state of the simulation.

The temporal argument, $t_{\mathbf{c}}$, is there because we should not consider the value of a cell as it is now, but rather the value that the cell would have by the time we get there (i.e., $t_{\mathbf{c}}$). If there is a very good fishing location inside the El Corralito area but I don't have time to get there before the seasonal closure, I should not value that location highly when choosing my destination. Note that $t_{\mathbf{c}}$ depends on the real travel time to the destination based on the speed of the fisher's vessel.

We write $\Omega$ as a short-cut for the general state of the simulation, but the aspects of that general state that are to be considered vary from function to function. For example, a vessel whose hold is almost full should give little weight to potential set opportunities, and more weight to the port. Regulations also come into play, whether they are temporal, spatial, or couched in terms of limits (e.g. on number of active FADs or sets).

All of those modulation functions are logistic functions of the same form as the scaling functions used for mapping the action values to the [0, 1] range (see Eq. 1) and the steepness and midpoints of those are likewise calibrated parameters of the model. If multiple functions apply to a field, their results are multiplied together, thus preserving the [0, 1] mapping. Table 3 provides a summary of the modulation functions used in the model.

**Table 3**   Modulation functions

| Field(s) | Factors influencing value |
| --- | --- |
| FAD, OFS, NOA, DEL | How full is the hold |
| | How close the fisher is to limit on number of sets |
| | Whether or not sets are allowed in that location at that time |
| OFS, NOA, DEL | How long since the location was last visited |
| DPL | How close the fisher is to the limit on active FADs |
| | Whether or not FAD deployments are allowed in that location at that time |
| PRT | How full is the hold |
| | How long it has been since departure |

## 10   Computing the Final Direction vector

In summary, we have: a set of attraction fields with values for each cell in the grid (Table 2), a set of modulation functions that allow to adjust these values according to circumstances (Table 3) and a set of weights (different for each fisher) that reflect their preferences for different types of actions (Fig. 3). All that is left do is is to combine all of those together in a single direction vector that tells the fisher which way to go.

In mathematical form, that final vector $\mathbf{d}$ is computed as:

$$g(a) = \sum_{\mathbf{c} \in C} \widehat{\mathbf{p_i} \mathbf{c}} \frac{\mathrm{G}^a_{\mathbf{c}_x \mathbf{c}_y} \prod_{\mu \in M_a} \mu(\mathbf{c}, t + t_{\mathbf{c}}, \Omega)}{(t_{\mathbf{c}})^n} \qquad (9)$$

$$\mathbf{d} = \sum_{a \in A} \widehat{g(a)} \cdot w_{i,a}, \qquad (10)$$

where, in addition to notation seen before, $g(a)$ is a function giving us the direction vector for action type $a$; $C$ is the set of all ocean cells; $\widehat{\mathbf{p_i} \mathbf{c}}$ is the unit vector (i.e., normalized so that it has a length of one) from the current position of the fisher ($\mathbf{p}_i$) to the ocean cell $\mathbf{c}$; $M_a$ the set of modulation functions that apply to field $a$; $t_c$, the time step at which the vessel would reach $c$ given its current location and speed; $t$ the current time step in the simulation; and $A$ is the set of action types {FAD, OFS, NOA, DEL, DPL, PRT}. Finally, the $n$ exponent used in $(t_{\mathbf{c}})^n$ is a calibrated parameter of the model that tells us how quickly the attraction of a cell declines with distance. Just as the attraction between two physical bodies is inversely proportional to the square of the distance between them, the attraction between a fisher and an ocean cell is inversely proportional to the $n^{th}$ power of the distance to that cell.

An example of the resulting vectors can be seen in Fig. 6. Note that agents in the model only compute the direction vector for the cell they are currently in, but looking at a larger portion of the ocean as in Fig. 6 allows us to see how following these vectors eventually leads a vessel to areas of high value.

The resulting vector $\mathbf{d}$ will almost certainly point to one of the eight neighbouring cells. In the improbable case where $|\mathbf{d}| = 0$, the agent stays put until circumstances change. When the agent arrives in the next cell, it will go back to applying the action strategy. If the next cell happens to be the port, the fisher unloads their hold, sells the catch, takes a few days off and then goes fishing again. And that concludes this very high-level overview of our gravity-based purse seine fishing behaviour algorithm.

**Fig. 6** Converting location values in attraction vectors

## 11   Empirical Calibration

As we walked through the model, we have identified various parameters that allow us to tweak how fishers behave, most notably the parameters of the logistic functions used in both the action and the destination strategies. To identify plausible values for these parameters, we have been working on calibrating the model against empirical data obtained from the IATTC observer database. We have identified various summary statistics that we believe to be the most important for judging if the model is a good representation of the fishery. We have focused on getting the number of actions of each type right across the whole fishery, and on reproducing the total catch for each species per type of action. In order to do this, we use a cluster-based niching evolutionary algorithm provided by the EvA2 library [15].

Figure 7 shows the targets used for calibration and the results from the best set of parameters that we have found so far. It's admittedly not great, but we are actively working on improving various part of the model (including but not limited to the behaviour algorithm) and we are hoping for better results shortly.

**Fig. 7** Preliminary calibration results. Black error bars represent standard deviation across runs

# 12    Conclusion

What we have presented here is just one layer in a much larger model, but it is the central one. By using a gravity model for this layer, we allow our agents to quickly adjust their behaviour in a very dynamic environment in which any long-term plan quickly becomes outdated. Our work is far from done, however.

As mentioned in the previous section, we are currently focussing on calibrating the model to empirical data. Once we are satisfied with the model's fit, we will shift our focus on comparing different policy scenarios, most notably alternative limits on the number of active FADs, limits on the number of sets that can be performed and individual catch limits per species.

We also intend to make two major improvements to the behaviour model. Our agents are already adaptive in the sense that they learn from experience about the value of different fishing locations, but we want to also give them the option to adapt their fishing strategies by adjusting their preferences for different actions according to how much profit each type of action generates for them.

Another crucial aspect that we haven't captured is coordination and communication within the fleet. We know that fishers in the same company often share FAD positions and that companies sometime make centralised decision as to who should fish on which FAD. Furthermore, we know that fishers across companies sometimes share information. FADs are sometimes given or traded. Locations of plentiful schools of fish are sometimes disclosed to friends or family working on other vessels. There is a rich world of agent interactions out there in the ocean, and we have yet to scratch its surface.

# References

1. Bailey, R.M., Carrella, E., Axtell, R., Burgess, M.G., Cabral, R.B., Drexler, M., Dorsett, C., Madsen, J.K., Merkl, A., Saul, S.: A computational approach to managing coupled human–environmental systems: the POSEIDON model of ocean fisheries. Sustainability Science **14**(2), 259–275 (2019). https://doi.org/10.1007/s11625-018-0579-9
2. Burgess, M.G., Carrella, E., Drexler, M., Axtell, R.L., Bailey, R.M., Watson, J.R., Cabral, R.B., Clemence, M., Costello, C., Dorsett, C.: Opportunities for agent-based modelling in human dimensions of fisheries. Fish Fisheries **21**(3), 570–587 (2020)
3. Carrella, E., Bailey, R.M., Madsen, J.K.: Repeated discrete choices in geographical agent based models with an application to fisheries. Environ. Modell. Softw. **111**, 204–230 (2019)
4. Carrella, E., Saul, S., Marshall, K., Burgess, M.G., Cabral, R.B., Bailey, R.M., Dorsett, C., Drexler, M., Madsen, J.K., Merkl, A.: Simple adaptive rules describe fishing behaviour better than perfect rationality in the US West Coast Groundfish fishery. Ecol. Econ. **169**(106), 449 (2020)
5. Chassignet, E., Hurlburt, H., Smedstad, O., Halliwell, G., Hogan, P., Wallcraft, A., Baraille, R., Bleck, R.: The HYCOM (hybrid coordinate ocean model) data assimilative system. J. Mar. Syst. **65**, 60–83 (2007). https://doi.org/10.1016/j.jmarsys.2005.09.016
6. Dagorn, L., Fréon, P.: Tropical tuna associated with floating objects: a simulation study of the meeting point hypothesis. Can. J. Fisheries Aquatic Sci. **56**(6), 984–993 (1999)

7. FAO (2020) The state of world fisheries and aquaculture 2020: sustainability in action. No: in the state of world fisheries and aquaculture (SOFIA). FAO, Rome, Italy (2020). https://doi.org/10.4060/ca9229en

8. Hall, M., Lennert-Cody, C.E., Garcia, M., Arenas, P.: Characteristics of floating objects and their attractiveness for tunas. In: Scott, M.D., Bayliff, W.H., Lennert-Cody, C.E. (eds.) Proceedings of the International Workshop on the Ecology and Fisheries for Tunas Associated with Floating Objects, February 11–13, 1992, Inter-American Tropical Tuna Commission Special Report 11, La Jolla, California (1992)

9. Lennert-Cody, C.E., Moreno, G., Restrepo, V., Román, M.H., Maunder, M.N.: Recent purse-seine FAD fishing strategies in the eastern Pacific Ocean: what is the appropriate number of FADs at sea? ICES J. Mar. Sci. **75**(5), 1748–1757 (2018). https://doi.org/10.1093/icesjms/fsy046

10. Lopez, J., Lennert-Cody, C., Maunder, M., Xu, H., Brodie, S., Jacox, M.J.H.: Developing alternative conservation measures for bigeye tuna in the Eastern Pacific Ocean: a dynamic ocean management approach. Tech. Rep. SAC-10 INF-D, Inter-American Tropical Tuna Commission, La Jolla, CA (2019)

11. Lopez, J., Román, M.H., Lennert-Cody, C.E., Maunder, M.N., Vogel, N., Fuller, L.M.: Floating object fishery indicators: a 2021 report. Tech. Rep. FAD-06-01 (2022)

12. Luke, S., Cioffi-Revilla, C., Panait, L., Sullivan, K., Balan, G.: Mason: a multiagent simulation environment. Simulation **81**(7), 517–527 (2005)

13. Madsen, J.K., Bailey, R., Carrella, E., Koralus, P.: From reactive towards anticipatory fishing agents. J. Simul. **15**(1–2), 23–37 (2021)

14. Orue, B., Lopez, J., Moreno, G., Santiago, J., Soto, M., Murua, H.: Aggregation process of drifting fish aggregating devices (DFADs) in the Western Indian Ocean: who arrives first, tuna or non-tuna species? Plos One **14**(1), e0210,435 (2019). https://doi.org/10.1371/journal.pone.0210435

15. Streichert, F., Ulmer, H.: JavaEvA—a java framework for evolutionary algorithms. Technical Report WSI-2005-06, Centre for Bioinformatics Tübingen, University of Tübingen, urn:nbn:de:bsz:21-opus-17022 (2005)

# Dynamics of Pedestrians' Flows During Daytime

**Marcin Wozniak** ⓘ

**Abstract**  The movement of people in the city varies significantly during the day. However, the availably of open localization data that could be useful in calibration of pedestrian ABM is negligible. The investigation of pedestrian traffic fluctuations could be an important element of city management (e.g. planning public transport, identification of bottlenecks). For that reason, the paper develops the agent-based model of pedestrians' flows dynamics in the center of one of the largest Polish cities (Poznan). The Google Places traffic data as well as census data and Geographical Information System were used to calibrate the model to generate reliable fluctuations of pedestrian movements. The developed ABM provides several valuable information that stand behind aggregate Google Places popular times rank. Mainly, we estimated the speed and size of pedestrians' flows together with the inflow and outflow of pedestrians to the city center. We were also able to identify bottlenecks, pedestrians' *waves* and areas of high/low density. The model captures and confirms several facts associated with fundamental diagrams of pedestrian flow and it could be used for further experiments regarding urban planning.

**Keywords**  Pedestrian traffic fluctuations · Geographic information system · Google Places

## 1  Introduction

Pedestrians' traffic vary during daytime. These fluctuations depend on the several external factors ranging from commuters [9] and tourists [12] to city inhabitants running their errands (e.g. [4]). The data on pedestrian behavior (including fluctuations of traffic) can be gathered through a variety of methods. These methods include field observations, experiments, survey methods, dedicated apps or localization data from mobile phones (e.g. [2]). Some of these methods are costly and time consuming and some of them require access to the highly vulnerable and sensitive

---

M. Wozniak (✉)
Faculty of Human Geography and Planning, Adam Mickiewicz University, Poznan, Poland
e-mail: woz@amu.edu.pl

corporation data.[1] Investigation of pedestrians' traffic dynamics may be important element of improving well-being and safety in cities (e.g. through adjustment of transportation network). Batty et al. [13] suggested that properly calibrated pedestrians' agent-based models could help both: in management of walkable environment and identify the potential (physical) barriers. The majority of agent-based models of pedestrian movements are strictly microscopic models [11]. These mainly focus on a details of pedestrians' interactions (e.g. transfer of emotions, panic in the crowd). Such approach improves the realism but also increases the computational demand and forecloses the analysis of more general and common scenarios. Therefore, in the paper we aim at these problems and develop pedestrian model set in the outdoor urban space which replicates fluctuations of pedestrian traffic in the city center. The model could be simply applied for any general scenario of daily traffic. We also propose the approach of calibrating the volume of pedestrian traffic which is based on easily available Google Places data. In turn, the outdoor environment and population density is calibrated and represented with Geographical Information System.

## 2 Data and Model

The simulation sandbox is the Old Market Square[2] in Poznan[3] with surrounding historic streets. The side of the area has approximately 450 m with the diagonal equal to 636 m. To obtain data on pedestrians' traffic we made use of simple traffic score provided by Google Places (GP). GP details provides several characteristics like ratings, addresses, opening hours or popular times (traffic data). The traffic data are displayed in a hourly bar chart and estimated from the localization data of mobile devices. We mapped GP scale with specific numerical values (Fig. 1). By assumption, the scale was set from 0 to 20 with 0.5 step (40 values). We chose Sunday as an illustrative example because the day covers all range of the traffic intensity values.



**Fig. 1** Pedestrian traffic on Sunday according to Google Places popular times (https://goo.gl/maps/MhXtUGaMUZ95ZnMVA) and adjusted numerical values

---

[1] See for example Rajpurohit et al. [8] for Facebook and WhatsUp case study.

[2] Brief description of the simulation area is available at: https://en.wikipedia.org/wiki/Pozna%C5%84_Old_Town.

[3] The city is situated in Western Poland. It is the fifth largest city in the country in terms of population (536,438 inhabitants) and sixth in terms of area (262 km$^2$).

Finally, we had to associate the numerical values for pedestrian traffic with specific numbers of pedestrians. We obtained data on the number of inhabitants of the Old Town quarter from the Web Map Service of the Poznan City Hall. According to census data, the Old Town quarter is inhabited by 1421 citizens in total. Therefore, we assumed that 1000 pedestrians can be at the same time on the Old Market Square as our simulation area covers 80% of the quarter. Therefore, the traffic scale was linked with specific number of pedestrians ranging from 0 to 1000.

During each simulation step, the agents are born at one of the simulation entry/ exit points according to the probability formula derived from GP data (Fig. 2). The concept is similar to the *gateways* described by Filomena et al. [3]. However, we did not associated these points with specific probability values. In the model, the gateways are situated in the locations associated with main streets, transport (tram, bus stops) as well as parking facilities. The general mechanics of agents' movement bases on gradient method described by Crooks et al. [1] with added enhancements regarding speed control module together with *discount* and *gravity* parameters that induce further heterogeneity in walking behavior (Fig. 2). The agents are seeking one destination point; the closer the agent is to the point, the stronger he/she is attracted by it. This point could be a shop, café or other popular place available at the Old Market Square. If agent finds destination, he or she exits the simulation by navigating to the one of randomly selected exit points. If agent reaches the point he or she permanently leaves the sandbox. In the model, the agents tend to travel at preferred walking speed = 1.4 m/s [14]. The simulation starts on Sunday at 4 a.m., ends on Monday 4 a.m. and covers 24 h, which are 54,000 ticks in the simulation sandbox (1 tick = 1.6 s). The time in the model was adjusted with the use of Sheppard's [10] NetLogo time extension. The general model logic together with benchmark parameters' values is presented in Fig. 2.

## 3   Simulation Results

The pedestrian traffic fluctuations produced by model covers the shape of Google Places bar plot (Fig. 1), however, due to advantages of ABM, one can have in-depth insight into this simple mechanism, e.g. the number of agents can be estimated. The maximum number of pedestrians was observed between 12 and 13 h (1008 pedestrians). The minimum values are observed in the early morning and at night.

The model reproduced some crowd inertia in the transitions periods between given hours. The inertia is represented through the outliers together with minimum/ maximum values. The pedestrian traffic does not switches instantly as the time evolves and needs some time to adjust to the new equilibrium conditions. The larger the change in traffic volume the adjustment time rises (Fig. 3).

Figure 4 plots density heatmaps and pedestrians trajectories. In the morning (8–9 a.m.) the pedestrians' flow is small and congestion zones are observed. Between

**Fig. 2** The general logic of pedestrians' dynamics model with calibration remarks. Model's parameters: **a** traffic parameter (dynamically adjusted on a basis of hourly Google Places data; vary from − 0.00065 to − 0.0023); **scent**—the value of destination patch (set to 200); **dist**—the distance from destination patch (in patches); **num_agents**—the initial number of agents in simulation (set to 5); **gravity**—diminishes the scent N(0.1, 0.1); **discount_rate**—diminishes the value of patches containing obstacles (set to 0.5); **min_speed**—minimum speed (set to 0.27 m/s); **max_speed**—maximum speed (set to 1.4 m/s)



**Fig. 3** The number of pedestrians in the simulation and pedestrian traffic according to Google Places

10–11 a.m. the traffic intensifies and some congestion zones are observed inside narrow historic streets. During peak hours (11–12), high density zones move from these streets to the areas where pedestrians try to leave the Old Market Square and enter one of the surrounding streets where entry/exit points are situated. This could be identified as kind of bottleneck scenario observed in the outdoor environment and

**Fig. 4** Density heatmaps and agents' trajectories during selected hours

reported by e.g. Luo et al. [6]. Figure 5 plots four main model outputs: the average speed and time agents spend in the simulation together with inflow of new agents and outflow of agents that leave the sandbox.

During daytime about 28,000 unique agents visited the Old Market Square. The highest inflow of new pedestrians was observed between 11 and 13 when it oscillated around 3000 agents per hour. In turn, the highest outflow of pedestrians (~2500) was noticed between 12 and 13. Such information could be extremely useful while planning transportation solution or other municipal services.

The average walking speed was 1.1 m/s (3.96 km/h) and the average time spend in the sandbox was 11 min (with standard deviation equal to 6 min). During the peak hours agents spend the longest time in the simulation (up to 25 min). The agents travel around preferred speed up to 10 a.m., then due to increased density they reduce pace by 20–30%. The lowest average speed is observed between 11 a.m. and 4 p.m. (~0.45 m/s) which are hours of the highest density. The standard deviation for speed was 0.31 m/s which corresponds with some previous findings regarding similar external conditions. E.g. Lam et al. [5] during empirical research identified mean speed in urban outdoor environment to be 1.19 m/s with standard deviation equal to 0.26 m/s; Older [7] observed mean speed and standard deviation in shopping streets is accordingly 1.3 and 0.3 m/s. Both of these pieces of research show the results that

**Fig. 5** Average speed, time spend in simulation (bars show standard deviation. The results are averages from 10 simulations.), and agents' flows (inflow and outflow)

are close to the output of developed model. We also found correlation coefficient between speed and density to be $-0.96$. The dependency matches fundamental diagrams of pedestrian flow (Campanella et al. 2009). Therefore, the model could be some important tool for exploring or planning pedestrians' solutions for improving the quality of city life.

# References

1. Crooks, A., Arie, C., Xu, L., Wise, S., Irvine, J.M., Stefanidis, A.: Walk this way: improving pedestrian agent-based models through scene activity analysis. ISPRS Int. J. Geo Inf. **4**(3), 1627–1656 (2015). https://doi.org/10.3390/ijgi4031627
2. Feng, Y., Dorine, D., Winnie, D., Serge, H.: Data collection methods for studying pedestrian behaviour: a systematic review. Build Environ **187**, 107329 (2021). https://doi.org/10.1016/j.buildenv.2020.107329
3. Filomena, G., Manley, E., Verstegen, J.A.: Perception of urban subdivisions in pedestrian movement simulation. PLoS ONE **15**(12), e0244099 (2020). https://doi.org/10.1371/journal.pone.0244099
4. Gorrini, A., Bandini, S., Vizzari, G.: Empirical investigation on pedestrian crowd dynamics and grouping. In: Chraibi, M., Boltes, M., Schadschneider, A., Seyfried, A. (eds.) Traffic and Granular Flow '13. Springer, Cham (2015). https://doi.org/10.1007/978-3-319-10629-8_10
5. Lam, W.H., Morrall, J.F., Ho, H.H.: Pedestrian flow characteristics in Hong Kong. Transp. Res. Rec. (1487), 56–62 (1995)

6. Luo, W., Jiao, P., Wang, Y.: Pedestrian arching mechanism at bottleneck in subway transit hub. Information **12**, 164 (2021). https://doi.org/10.3390/info12040164
7. Older, S.J.: Movement of pedestrians on footways in shopping streets. Traffic Eng. Control **10**(4), 160–163 (1968)
8. Rajpurohit, D., Singh, G., Yadav, K.: A socio-legal analysis of whatsapp privacy policy 2021 in India. Contemp. Stud. https://doi.org/10.2139/ssrn.3850579
9. Ren, H., Song, Y., Li, S., et al.: Two-step optimization of urban rail transit marshalling and real-time station control at a comprehensive transportation hub. Urban Rail Transit. **7**, 257–268 (2021). https://doi.org/10.1007/s40864-021-00157-4
10. Sheppard, C.: NetLogo Time Extension (2017). https://github.com/colinsheppard/time
11. Tordeux, A., Lämmel, G., Hänseler, F.S., Steffen, B.: A mesoscopic model for large-scale simulation of pedestrian dynamics. Transp. Res. Part C Emerg. Technol. **93**, 128–147 (2018)
12. Zhang, B., Li, N., Shi, F., Law, R.: A deep learning approach for daily tourist flow forecasting with consumer search data. Asia Pacific J. Tour. Res. **25**(3), 323–339 (2020). https://doi.org/10.1080/10941665.2019.1709876
13. Batty, M., Jiang, B., Thurstain-Goodwin, M.: Local movement: agent-based models of pedestrian flows. (CASA Working Papers 4, 1998). Centre for Advanced Spatial Analysis (UCL): London, UK
14. Betty J., Mohler William B., Thompson Sarah H., Creem-Regehr Herbert L., Pick William H., Warren.: Visual flow influences gait transition speed and preferred walking speed Experimental Brain Research **181**(2), 221–228 (2007). https://doi.org/10.1007/s00221-007-0917-0
15. Campanella M, Hoogendoorn SP, Daamen W.: Effects of Heterogeneity on Self-Organized Pedestrian Flows Transportation Research Record: Journal of the Transportation Research Board **2124**(1), 148–156 (2009). https://doi.org/10.3141/2124-14

# Evacuation Simulation for Large-Scale Urban Population

**Etzion Harari, Naphtali Abudarham, and Tomer Rokita**

**Abstract**  Large-scale population evacuation from urban areas may occur during disasters such as earth quakes, volcano eruptions, militarized conflicts, environmental disasters and more. Efficient and safe population evacuation is of great importance as it can save lives and reduce human suffering. The current study demonstrates an Agent-Based Simulation tool which may be used to support operational planning for population evacuation from threatened urban areas. The simulation models households as agents, each acting in accordance to a designated decision function, which renders the probability of evacuation as a function of the socioeconomic and demographic characteristics of the agents and the behavior of neighboring ones. Upon evacuation decision, agents embark on their way to their preassigned destinations, while their optimal route is calculated and updated periodically, based on road information (taken from Open-Street Map), accumulative traffic congestion, and simulated road conditions. The simulation calculates and records the location of all agents and enables the user to identify and analyze different evacuation scenarios, compare evacuation sequences, map and identify road bottlenecks, etc. Integrating such a simulation into the planning process—both at the municipal and the national levels—can significantly enhance authorities' processes of preparing evacuation plans, including investing resources for developing safe evacuation destinations and educating the population for the unfolding of future emergencies. The main contribution of the current work is the ability to efficiently calculate optimal routes for millions of agents. To demonstrate this we applied the simulation on Kyiv (Ukraine), where a large number of its 3 million citizens have fled the city during Russia's invasion on February 2022.

**Keywords**  Agent-based simulation · Population evacuation · Disaster areas · Routing algorithms

E. Harari · N. Abudarham · T. Rokita (✉)
Rafael—Advanced Defense Systems Ltd., Gazit Institute, Tel-Aviv, Israel
e-mail: rokita.tomer@gmail.com

# 1 Introduction

Environmental disasters and militarized emergencies are not unique in today's international landscape. Such emergencies render considerable human suffering for those who find themselves detached from their safe homes in quest for shelter. Despite the recurrent nature of such events, authorities as well as the common people are not adequately prepared for large-scale evacuation, as most of them are required to take action which is not always optimal in terms of efficiency, or safety. This is largely due to the poor preparedness of both the public and the authorities in terms of assessing how long the evacuation process takes, where can the population find shelter and so on. The uncertainty accompanied to much of contemporary populations' evacuation may result in heavy traffic congestion, over-crowding of shelters and other undesirable progressions, which, in turn, may increase tension, or obstruct the attempt to attain safety and security. In extreme cases, unplanned evacuations may lead to heavy casualties. Emergencies that require large-scale population evacuations may include earth quakes, volcano eruptions, disasters such as leakage of dangerous chemicals from plants, or even wars. In all of these cases operational planning for efficient and safe population evacuation is of considerable importance.

The behavioral proclivities of populations facing eminent threats or emergencies have traditionally constituted one of the most complex and challenging research fields in social sciences in general, and in the study of human behavior in particular [9, 10, 12, 13, 17, 19, 20]. In recent decades, emergency management leaders at the state and municipal levels increasingly stress the need for effective disaster preparedness strategies, which enable effective evacuation of populations from threatened areas [14].

To plan effective evacuation two major components are taken into account: the decision of individuals whether or not to evacuate and to what destination, and the dynamics of the motion of the population (i.e. the traffic). The majority of the scholarly literature on evacuation behavior argues that population flows during disasters are not the aggregation of individuals' decisions, but rather the result of a complex interactions of various factors at different levels (mostly the individual and the environmental/societal) [8]. Chiefly among them are the following:

1. Early Warning and Signal Systems: the main means through which the individual interprets reality. These include the media (TV, Radio, Smartphones) and more "traditional" vectors such as alarms, heralds, etc. The extent of emergency and clarity attributed to the information that is being mediated to the population via such vectors are of profound importance for passing and internalizing the evacuation messages throughout the population [1].

2. Risk Perception: the vast majority of the scholarly literature on evacuation behaviors considers risk perception as the variable which holds the greatest explanatory power for populations' behavioral proclivities in emergencies. This variable is accumulative in nature as it constitutes a synthesis of several explanatory factors such as : past experience (has the population experienced a similar emergency in the past and how did it react at the time?); Socio-demographic variables (such

as marital status, young children at the household); Vulnerability (vulnerable populations); Car ownership (a significant factor in one's decision to leave independently of others in the threatened area); Infrastructural factors such as the physical/geographical proximity of populations to the threatened/damaged area and the quality of the infrastructure surrounding them (roads, hospitals, shelters, public transportation); and lastly, psychological factors as the extent to which the population perceives the seriousness of the threat and its potential infiltration upon them.

3. Social Resources: overall, people who evacuate out of their homes will have a significantly larger number of (closer) social connections, than those who choose to stay put [3, 16]. Relatives and family members (immediate and/or extended), friends and acquaintances residing in "safe zones" are considered, then, as accelerating populations to evacuate. Evacuation destinations are related to the geographic distribution of one's social network, i.e. people tend to move to places where their friends or family are located [4, 7].

4. Environmental Resources: Facilitating evacuation infrastructures:

(a) Public Infrastructures: shelters, schools, public halls, stadiums, soccer fields and any large and wide facility—far enough from the disaster zone—that can host evacuating populations and provide them with physical security (even at the most minimal level).

(b) Social-Communal Infrastructures: open spaces which enable evacuating populations a brief stay (religious institutions, non-governmental organizations' properties etc.).

These factors further reinforce the operational planning of population evacuation as a complex and highly delicate task; such that requires mass exploration, identification, extraction, processing and analysis of data on the relevant population and the environment within which it resides and with which it interacts in order to better plan and prepare the potential time, the space, the destinations for populations' evacuation behavior, as well as to better prepare and assess the potential ramifications and implications of the differentiated behaviors of populations.

To this end, in the current study an Agent-Based Simulation is developed, which can incorporate some of the key factors described herein, and serve as a planning tool which enables authorities to draw more effective and safer evacuation plans for different disaster or emergencies scenarios, as well as mapping and analyzing the environment under their jurisdiction, identify bottlenecks, prepare solutions such as clear and known evacuation destinations and shelters, safe roads and public transportation, as well as guidelines and population education for future occurrences. The choice of Agent-Based Modelling (**ABM**) method relies on the acknowledgment that it is the most effective tool for modelling complex phenomena whose dynamics are determined by the interactions of multiple elements. A common example for the use of ABM is traffic analysis: while each driver may act based on simple individual rules, traffic jams cannot be analytically predicted, they can only be simulated.

ABM has been widely used to simulate evacuation in various scenarios such as buildings in case of fire [11, 18], evacuation from sport stadiums [22] and tzunami

events [2, 21]. However, to the best of our knowledge, no previous study has simulated the evacuation of whole cities, taking into account the motion of millions of people. Obviously, as the size of the simulation increases so does the difficulty of accurately modelling the behavior of such large numbers of agents. However, just like any large scale simulation, the purpose of using our tool is not to precisely predict the action of each agent, but to enable gross comparison of different strategies and to assist the authorities in gaining insights which may be difficult to achieve by static analysis.

## 2 Evacuation Simulation

### 2.1 General Description

The simulation contains two main components: Households (families) as agents, and roads and evacuation destinations. To drive the simulation we use evacuation triggers: each trigger has a type (e.g. evacuation notification), geographic area of influence (**AOI**), and the time in which it occurs. Once a trigger occurs, an evacuation decision function of each agent within its AOI is evaluated, to determine whether the agent should begin evacuation. The evacuation decision function may be defined by the user; in the current study we defined a function that takes into account socioeconomic status (high or low) and demographic data (does the family include young children), as well as environmental considerations (how many agents in the neighborhood have already evacuated), as described in Sect. 3.1. Once an agent decides to evacuate, the optimal route is efficiently calculated to its preassigned destination, based on the current traffic and road status. At every time step the simulation advances all moving agents along their routes, at a speed that is determined by the current traffic. In addition, every predetermined time interval the simulation calculates new optimal routes for all moving agents. Road conditions may change during the simulation (e.g. roads destroyed by bombings, roads blocked by police), requiring agent route recalculation. This approach of enabling agents to calculate optimal routes based on full information on the roads and other agents is, of course, unrealistic. However, this approach is useful in focusing the attention of authorities/planners to where bottlenecks appear even if all agents take optimal routes. The simulation records the positions of all agents throughout the simulation (which runs for a predetermined number of time steps), enabling the user to analyze the data in various ways, for example the evacuation status at each time step, number of evacuated agents by time, traffic jams and much more. Finally, the user may run the simulation many times, altering the data such as evacuation triggers (type, time and location), roads data (evacuation destinations, road blocking), agent attributes, agent decision function, etc., and compare the different outcomes in order to derive the best evacuation plan for different scenarios.

## 2.2  Households as Agents

An agent in the current model is a household (or a family), which acts as a single unit that decides whether or not to evacuate its home and head to an evacuation destination (e.g. a shelter, a city exit point etc.). The agents are initially located on the city map based on the user's input, which can be taken from official municipal databases. All agents begin in status "dormant", and move to status "active" if the output of the invocation of their evacuation decision function was that they should evacuate. Agents remain in status "active" until they reach their evacuation destination, where their status changes to "saved". Not all agents who became active end up is the "saved" status, for example if they got stuck in traffic or in blocked roads.

## 2.3  Roads and Evacuation Destinations

Information on the roads and evacuation destinations are taken from the Open-Street-Map (**OSM**) database [15]. For example, Fig. 2 shows the the full road map of Kyiv. We convert the road map to a directed graph, where each road junction (a point connecting at least three roads) is a node (the set of all nodes is denoted as $V$) and the roads connecting junctions are graph edges ($E$). If a road is bi-directional it will be entered twice to the graph, once for each direction. For the Kyiv test-case, the road graph contains 80,000 nodes and 200,000 edges.

A graph is defined using the nodes $V$ and the adjacency matrix $A$, where the value in cell $A_{i,j}$ is 1 if $V_i$ and $V_j$ are connected, and 0 otherwise.

$$G = (V, A) \tag{1}$$

## 2.4  Efficient Optimal Route Calculation

An optimal $P_{ij}$ route is calculated for each evacuating agent from its point of origin ($V_i$) to its designated evacuation destination ($V_j$). This calculation is performed for every agent when it starts its evacuation, and again when road conditions change and at predetermined intervals. Normally we would use some efficient shortest path algorithm (such as $A*$) to calculate paths from source to destination, however, in the current work we need to calculate routes for millions of agents, so we took a different approach, using the Dijkstra shortest path algorithm [5]. As an initial step, Dijkstra's algorithm searches the entire graph, starting from the source point, and constructs a data-structure $P_i$. which holds the shortest paths from the source point to all other points in the graph. Instead of calculating the paths from agent locations to evacuation destinations, we calculate the reverse paths—from evacuation destinations to the rest of the graph: Fig. 1a shows the initial graph $G$, with directed edges, and destination

a) Given Directed Graph $G$

b) Reversed Graph $G^t$

c) Dijkstra Shortest Paths from D

d) Paths from every node to D



**Fig. 1** Calculating the shortest paths to the destination point ($D$) using $R$-Dijkstra

node D. Then in Fig. 1b we create the reverse graph $G^t = (V, A^t)$. Then, in Fig. 1c we calculate $P_{D.}$ the shortest paths from the destination $D$ to all other nodes. Finally, in Fig. 1d we reverse the edge directions again, and immediately obtain $P_{.D}$ the shortest path from every agent (i.e. from any node in the graph) to the evacuation destination. Since the number of evacuation destinations is very small (less than 10), we end up running this calculation a very small number of times to get the paths for millions of agents. For example, calculating the paths for 1,000,000 agents in our test case using A* takes **75 min**. However, calculating the same paths using $R$-Dijkstra algorithm takes just **3.5 s**, which is 1250 times faster then the naive $A*$ method.

The Dijkstra implementation in NetworkX [6] Python package, allows us to build a navigation engine that uses the Reversed-Dijkstra algorithm. This code uses the road graph, to calculate the path to each of our $T$ destination points. In addition, similar evacuation destinations can be grouped together (e.g. all the subway stations) by adding "logical" evacuation destinations (as demonstrated in Fig. 2), further reducing the size of $T$.

For calculating the shortest paths as well as traffic congestion the length and width of each road was obtained from OSM, and these values were assigned to the graph edges. The capacity of each edge is calculated by multiplying the edge length with the road's width, and dividing by a constant value representing the "volume" of an agent. At each step of the simulation all agents advance on their path at a constant speed, so the travel time is a linear function of the edge length. The simulation records the number of agents at each edge, and if the number of agents on an edge exceeds its capacity, the length of the edge is artificially increased, which results in longer travel time, and of course affects the shortest path calculation.

**Fig. 2** The map of Kyiv with division to city regions (left pane), and the road information taken from OSM (right pane). Blue circles indicate subway stations, which are connected to a single logical node

## 3   Case Study—Kyiv, Ukraine

In the following section the application of the simulation is demonstrated on the city of Kyiv, in the context of the Russian invasion of February 2022. Different evacuation scenarios were tested, demonstrating the types of insights that are possible to obtain using the simulation.

### 3.1   Kyiv Simulation Data

1. **Kyiv population data**: The number of people within each city region was taken from Wikipedia, reaching a total of 2.75 million people.[1]
2. **Simulation agents**: The average household size in Ukraine is 2.58, therefore 1,060,000 agents were created. The agents were randomly positioned within the city regions based on region populations. To support the agents' evacuation decision function each agent was assigned the following attributes:

    (a) Socioeconomic status: High or Low. Due to lack of official data this value was determined based on the density of buildings in the agent's neighborhood, assuming that low building density indicates a high socioeconomic status, and high density—low socioeconomic status. Overall 50% of the agents were

---

[1] Kyiv region was obtained map from: https://geodata.lib.utexas.edu/catalog/stanford-pp624tm0074.

assigned high socioeconomic status, and their distribution across the city was according to building density.

(b) Does the household include young children? Yes or No. Due to lack of official data this value was estimated based on agent distance to schools or playgrounds: a small distance indicated a family with small children. Overall 10% of the agents had young children.

3. **Agent evacuation decision function**: the following evacuation decision function was designed, based on agent attributes, displayed here as a decision tree:



The percentages at the leaves indicate the probability of the agent to evacuate ($EvacP$) its home when its decision function is activated, following an evacuation trigger in its location. This probability is then adjusted according to the percent of already evacuated agents the current agent's 100 meter radius($EvacRate$). The final evacuation probability is then adjusted as follows:

$$Final\ Evacuation\ Probability = EvacP + (1 - EvacP) * EvacRate$$
(2)

4. **Kyiv road map and evacuation destinations**: Kyiv road map was obtained from OSM. Evacuation destinations were determined to be the western exits from the city and the subway stations (from OSM transportation layer). The preassigned evacuation destination of 80% of the agents with high socioeconomic status was the western city exits, and the rest of them were directed to subway stations. The reverse was applied to the agents with low socioeconomic status. Figure 2 shows the map of Kyiv with division to city regions (left pane), and the road information taken from OSM (right pane). The locations of the subway stations are also shown, which are all connected to a single logical node, which is used to accelerate the computations (as explained in Sect. 2.4).

## 3.2 Analysis of Evacuation Scenarios

1. **Traffic bottlenecks identification**. The purpose of this scenario was to identify road bottlenecks in Kyiv in case of a large-scale evacuation. To this end we generated evacuation triggers in the form of notifications to the public. All city regions received 5 evacuation notifications simultaneously during the first 10 minutes of the simulation (one notification every 2 minutes), thereby activating

the decision functions of all the agents in the city. A total of 700,000 agents embarked on their way (representing 2,000,000 citizens), and the simulation was ran for 1000 time-steps (each time-step simulated 1 min, making it a total of approximately 17 simulated hours) The same simulation was ran again, with blocked roads heading north-west out of Kyiv, to simulate Russian attacks (some of which came from the direction of Belarous). Figure 3 shows the road status at the end of the simulation. The left image shows road congestion—the red roads are the critical bottlenecks. The right image shows the same scenario when the north-western roads are blocked, and it can be seen that the critical roads (bottlenecks) are different, almost eliminating the west road bottlenecks.

Figure 4 shows the effects of the evacuation decision function. The left pane shows the percent of evacuated agents by socioeconomic status. It can be seen for example that after 200 minutes 60% of the population classified as low socioeconomic status has reached their evacuation destinations, and this is in accordance with the fact that 80% of this population had subway stations as their preassigned evacuation destinations. In contrast, it took nearly 1100 minutes to evacuate the high socioeconomic populations, which mostly headed to the city exits. This kind of analysis demonstrates the importance of planning optimal routing (as well as considering the capacity of in-city shelters like subway stations, which is not taken into account in the current study). The right pane shows the percent of agents that were activated (i.e. decided to evacuate their homes) following the evacuation triggers. These numbers are correlated with the evacuation decision function described above.

2. **Evacuation of adjacent versus nonadjacent areas**. In this experiment, simultaneous evacuation of a similar number of agents, that are either located in adjacent or nonadjacent areas, is simulated. As shown in Fig. 5, when agents in adjacent areas attempt to evacuate the city simultaneously, the evacuation rate is lower than the case of evacuating agents from non-adjacent areas. This is because the traffic congestion is larger when a large number of agents move out of a single area than when agents move out of scattered areas. This result is crucial if one wants to minimize evacuation time, as it directs authorities to carefully plan evacuation sequence in order to avoid unnecessary delays. This result is not trivial; one could not easily reach this kind of insight without a tool such as this simulation.

## 4   Conclusions and Future Work

The current study demonstrates an ABM simulation model for large-scale evacuation of an urban population. The simulation incorporates the following elements:

1. Agent behavior model: using a decision function which determines whether the agent will evacuate its home in response to an environmental signal, based on the

**Fig. 3** Road congestion in Kyiv area. The left image shows road congestion, the red roads are the critical bottlenecks. The right image shows the same scenario when the north-western roads are blocked, to simulate Russian attacks



**Fig. 4** The effect of the agent evacuation decision function. Left pane: percent of evacuated agents by socioeconomic status. Right pane: percent of agents that were activated following the evacuation triggers

     agent's attributes, such as socioeconomic and demographic status, the behavior of neighboring agents and more;
2. Efficient large-scale road traffic model based on OSM data.

Taken together, the developed simulation can be useful for authorities to test and compare evacuation scenarios, identify road bottlenecks, prepare infrastructure to support safe and efficient evacuation, educate and prepare the population, and much more. In this paper a single case-study was described, for which very little official information about agent attributes, such as socioeconomic status and demographics, was available. In real-life scenarios, authorities could obtain much more accurate information, and use tools such as public surveys to study the actual stands and predispositions of the population, to create a much more accurate evacuation decision

**Fig. 5** Evacuation rate for evacuating agents from adjacent (red rectangle) versus non-adjacent (blue squares) area. The total number of agents in the red and blue areas is the same

function, to assign preferred evacuation destinations, and more. The authors hope that this work can inspire authorities to adopt this approach and use such tools to better prepare for unfortunate events.

# References

1. Chadefaux, T.: Early warning signals for war in the news. J. Peace Res. **51**(1), 5–18 (2014)
2. Chen, C., Koll, C., Wang, H., Lindell, M.: An interdisciplinary agent-based evacuation model: integrating natural environment, built environment, and social system for community preparedness and resilience. Nat. Hazards Earth Syst. Sci. Disc. **2021**, 1–27 (2021)
3. Collins, J., Ersing, R.L., Polen, A., Saunders, M.: Evacuation behavior measured during an evacuation order: an assessment of the effects of social connections on the decision to evacuate. Nat. Hazards Center (2018)
4. Dash, N., Gladwin, H.: Evacuation decision making and behavioral responses: individual and household. Nat. Hazards Rev. **8**(3), 69–77 (2007)
5. Fredman, M.L., Tarjan, R.E.: Fibonacci heaps and their uses in improved network optimization algorithms. J. ACM (JACM) **34**(3), 596–615 (1987)
6. Hagberg, A., Swart, P., S Chult, D.: Exploring network structure, dynamics, and function using networkx. Tech. rep., Los Alamos National Lab.(LANL), Los Alamos, NM (United States) (2008)
7. Hasan, S., Ukkusuri, S., Gladwin, H., Murray-Tuite, P.: Behavioral model to understand household-level hurricane evacuation decision making. J. Transp. Eng. **137**(5), 341–348 (2011)
8. Hong, L., Frias-Martinez, V.: Modeling and predicting evacuation flows during hurricane IRMA. EPJ Data Sci. **9**(1), 29 (2020)
9. Houts, P.S., Cleary, P.D.: Three Mile Island crisis: psychological, social, and economic impacts on the surrounding population. Penn State Press (2010)
10. Iklé, F.C., Kincaid, H.V.: Social Aspects of Wartime Evacuation of American Cities: With Particular Emphasis on Long-Term Housing and Reemployment, vol. 4, National Academy of Sciences, National Research Council (1956)
11. Kasereka, S., Kasoro, N., Kyamakya, K., Doungmo Goufo, E.F., Chokki, A.P., Yengo, M.V.: Agent-based modelling and simulation for evacuation of people from a building in case of fire.

Procedia Comput. Sci. **130**, 10–17 (2018). https://doi.org/10.1016/j.procs.2018.04.006, https://www.sciencedirect.com/science/article/pii/S1877050918303569, The 9th International Conference on Ambient Systems, Networks and Technologies (ANT 2018) / The 8th International Conference on Sustainable Energy Information Technology (SEIT-2018) / Affiliated Workshops

12. Kuligowski, E.: Predicting human behavior during fires. Fire Technol. **49**(1), 101–120 (2013)
13. Newman, S.M.: The occasion instant: the structure of social responses to unanticipated air raid warnings (1963)
14. Norris, F.H.: Disaster research methods: past progress and future directions. J. Traumatic Stress: Off. Publ. Int. Soc. Traumatic Stress Stud. **19**(2), 173–184 (2006)
15. OpenStreetMap contributors: Planet dump retrieved from https://planet.osm.org, https://www.openstreetmap.org (2017)
16. Paton, D., Johnston, D.: Disaster resilience: an integrated approach. Charles C Thomas Publisher (2017)
17. Raphael, B.: Individual and Community Responses to Trauma and Disaster: The Structure of Human Chaos. Cambridge University Press (1995)
18. Ren, C., Yang, C., Jin, S.: Agent-based modeling and simulation on emergency evacuation. In: Zhou, J. (ed.) Complex Scieces, pp. 1451–1461. Springer, Berlin Heidelberg, Berlin, Heidelberg (2009)
19. Silver, R.C., Holman, E.A., McIntosh, D.N., Poulin, M., Gil-Rivas, V.: Nationwide longitudinal study of psychological responses to September 11. JAMA **288**(10), 1235–1244 (2002)
20. Tierney, K.: Disaster beliefs and institutional interests: recycling disaster myths in the aftermath of 9–11. In: Terrorism and Disaster: New Threats, New Ideas. Emerald Group Publishing Limited (2003)
21. Zhan, F.B., Chen, X.: Agent-based modeling and evacuation planning. In: Geospatial Technologies and Homeland Security, pp. 189–208. Springer (2008)
22. Zhang, W., Terrier, V., Fei, X., Markov, A., Duncan, S., Balchanos, M., Sung, W., Mavris, D., Loper, M., Whitaker, E., Riley, M.: Agent-Based Modeling of a Stadium Evacuation in a Smart City. pp. 2803–2814 (12 2018). https://doi.org/10.1109/WSC.2018.8632176

# Extending Partial-Order Planning to Account for Norms in Agent Behavior

**Tokimahery Ramarozaka** ⓘ**, Jean-Pierre Müller** ⓘ**,
and Hasina Lalaina Rakotonirainy** ⓘ

**Abstract**  Following a couple of models aiming to assess the effectiveness of norms in Madagascar on the MIMOSA platform, Müller et al., have noticed that the current architecture was not expressive enough to deal with all relevant norms, their different aspects, and how they interfere with the agent's behavior for such complex systems. In response, this paper proposes a new agent architecture and its dedicated language to enhance the expressiveness of norms in agent-based modeling. The architecture has to (1) identify all the applicable norms given a temporal, spatial, and social context, and (2) generate an agent behavior to account for these norms. We propose to extend automated planning and use a Model-Driven Engineering approach to build the abstract and concrete syntaxes of the language and its semantics. The resulting architecture will allow modelers to express a wider spectrum of norms and provides a normative decision tool that will ease further discussions and interpretations.

## 1  Introduction

To answer the questions of norm effectiveness on renewable resource management, agent-based models (ABM) were explored with the MIRANA [1] and HINA [2] models, through the MIMOSA simulation platform [3]. However, both model's implementations were made ad-hoc, without any generic structure of norm, nor how they affect agent behavior, making them hard to (re)use and understand. Moreover, their lack of specification on a norm's spatial (where is a norm applicable?), temporal (when is it applicable?), and social context (to whom do they apply?) does not allow

T. Ramarozaka (✉) · H. L. Rakotonirainy
Informatique, Géomatique, Mathématiques et Décisions (IGMA), Andrainjato, University of Fianarantsoa, Fianarantsoa, Madagascar
e-mail: tokyramarozaka@gmail.com

J.-P. Müller
CIRAD—UMR SENS, Campus International de Baillarguet, 34398 Montpellier, France

modelers to express all relevant norms for such complex topics. For example, a norm that states that "it is forbidden to fish in winter around mangroves." might cause an agent to wait for the next season to fish, or to fish in non-mangrove areas; but how an agent reason and generate such behavior autonomously remains elusive.

While the literature on norms in ABM clearly offers a handful of normative agent architectures, such as BOID [4], NoA [5], or N-BDI [6], most of them focus on extending the BDI (Belief Desire Intention) architecture [7] with norms and does not provide the needed spatial, temporal and social context to be more expressive about norms. Recent works in [8] offer a mean to express these contexts and compute all applicable norms in that regard, but without describing how they are being applied in the agent decision-making process, i.e., in the agent architecture.

In response, this paper aims to propose a new agent architecture with its dedicated agent architecture to account for norms in agent behaviors while accounting for the spatial, temporal, and social context of norms. By using Model Driven Engineering (MDE) to build the language, our goal is to provide a tool to test the effectiveness of different norms on the agent's behavior and the whole socio-ecosystem.

Section 2 presents and justifies the methods we used to describe norms, and their influence on the agent's behavior. Section 3 then describes the results. Section 4 discusses how relevant they are for the issue at hand. Section 5 concludes by giving an overview of the current state of the project, and perspectives on future works.

## 2 Related Works

### 2.1 Norm Definition

In ABM, a norm is generally defined as a soft constraint on the agent's behavior [9, 10] to attain a desired behavior in time and space [11]. We represent this regulating mechanism through institutions [12], which is a set of ontology and norms. Each institution is endowed with a number of norms about the various social roles for the agents.

For instance, an institution could be a village where the *fisherman* role ought to fish only in authorized areas and have a license, while the *chief* role ought to provide enough food for his village. Norms can either tell what ought to be done under certain conditions, i.e. regulative norms; or, define that something counts as something else for a given institution, i.e. constitutive norms.

Constitutive norms introduce abstract classifications that some existing objects count as some concept or role within the institutions. They allow us to instantiate institutions as organizations by stating that an individual *John* counts as a fisherman, or that *Peter* counts as *chief* and that a certain place counts as a village.

Regulative norms are generally divided into three categories:

- Obligations: some action that the agent ought to do, or some fact that an agent ought to make true in a given situation;
- Prohibitions: some action that the agent ought not to do, or some fact that should not be true in a given situation;
- Permissions: some action or fact that is permitted in a certain situation(s) [13].

We chose the ADICO (Attribute Deontic aIm Condition Otherwise) formalism to express norms [14] as it covers most of the elements of norms in MAS, such as social roles describing the desired behaviors, rewards, and sanctions when applying or violating them [15], and more recently the notions of time and space formalized in [8, 16].

## 2.2 Model of Behavior

In classical AI, automated planning [17] is often used to generate agent behavior. It needs to define a planning problem with: some initial situation, a goal, and a set of possible actions, and outputs a sequence (or a partially ordered set) of actions called the *plan*, which allow the agent to reach its goal. The planning problem is then assimilated to a state-space search: a set of possible states is explored using operators to compute the next states until we find a solution state. Depending on the chosen approach, states and operators may encapsulate different concepts.

Since performance is not the main focus here, we chose Partial-Order Planning (POP) amongst other approaches for its flexibility: (1) states in POP are partial-order plans, they set constraints on the arrangement of the plan, rather than set a fixed sequence which would be harder to adapt when new or unexpected changes occur in the environment; (2) it is the most suited for agent planning as advocated in [18] because of the complexity of agent-based models. Each newly generated state (= plan) is examined for possible flaws like unachieved goals or conflicting actions. We compute the next state by solving one of those flaws, either: adding in new actions, arranging the partial-order of the plan, or constraining the variable bindings of its existing actions.

## 3 Our Contributions

### 3.1 Extending Planning Problems

A planning problem is defined by: some initial situation, a goal, and a set of possible actions, as summarized by Fig. 1. We represent actions in a deterministic way as in STRIPS [19], i.e. the consequences of the action are certain and occur simultaneously

**Fig. 1** Structure of a classic planning problem

when the action is executed. Dealing with conditional consequences is beyond the scope of this paper.

We propose to extend the planning problem (Fig. 1) with the concept of norms through organizations and institutions. Norms are applied to an agent through its roles within organizations, which are instances of one or more institutions.

For example, if the *fisherman* role is prohibited to fish in the *village* institution without a license, and the agent *Jon* plays the role of a *fisherman* in some *village v1* (*Jon* counts as a *fisherman*, *v1* counts as a *village*), then all the norms which would apply to *fisherman* apply to *Jon*: he is prohibited from fishing in *v*1 without a license. Thus, to define a planning problem with norms, we must describe the organizations in which the agent plays a role(s), listing all applicable norms, as shown in Fig. 2.

Each norm has a deontic operator (obligation, prohibition, permission) and applicability conditions with consequences that describe which action or proposition is mandatory, prohibited or permitted.

## 3.2   Extending POP's Internal Structure

States in classical POP are defined as plans, which we improve by correcting flaws to obtain an executable plan achieving the goal (if any). Formally, Fig. 3 defines a plan as a tuple <*S, A, Cc, Tc, F*> where:

- *S* is a set of situations;
- *A* is a set of steps, i.e. instances of actions;
- *Cc* is a set of (non) codenotation constraints describing that some variable must be bound to some value or not;
- *Tc* is a set of temporal constraints that describes which situation/step must be before another, defining a partial order;
- F is a set of flaws, which can either be open conditions: preconditions of steps that need to be satisfied, or threats: two conflicting steps, as stated in [20].

To compute the next state, a flaw in the plan is chosen and resolved.

**Fig. 2** Our extensions (highlighted in red) of a planning problem with norms



**Fig. 3** Internal structure of a state in POP

Plan modification simply consists in adding situations, steps, (non) codenotation and/or temporal constraints to resolve a flaw. The planning stops when it has found an executable plan satisfying the goal. A plan is executable if, for each of its steps, all of its preconditions are necessarily true in its preceding situation, i.e. it does not contain open conditions; and its *Tc* and *Cc* do not contain contradictions. By introducing a dummy final step with the goal as its precondition, this executability condition is enough to find a solution plan (if any).

**Fig. 4** Our extensions (highlighted in red) on the internal structure of a state in POP

Each state or partially ordered plan can be instantiated as a set of complete plans. A plan is complete if its $Tc$ (temporal constraints) define a total order on its situations and steps, and its $Cc$ (codenotation constraints) grounds all variables to a constant.

Throughout this paper, we shall refer to a proposition being «necessarily true», following Chapman's modal truth criterion [20] if that proposition is satisfied in any complete plan specified by a (partial) plan. We extend this existing structure with the notions of norm and interval, as highlighted in red in Fig. 4.

Since we introduce new flaws related to norms, we need to describe which plan modifications (operators) can resolve them.

### 3.3 Dealing with Norms

Our extensions to allow agents to deal with norms in the decision-making process rely on three key strategies: respect, violation, and circumvention. In the following section, we will take a look at how agents can respect or circumvent norms, and then see how they can violate them.

**Dealing with obligations**. Obligations are mandatory actions or propositions that need to be necessarily true in a certain situation of the plan. To consider them, we define a new flaw in POP.

*Definition 1 (missing obligation)*. A missing obligation is a flaw such that the applicability conditions of an obligation are necessarily true in a situation $S$, and:

- the mandatory action $A$ is not succeeding to $S$;
- or, the mandatory proposition $P$ is not necessarily true in $S$

Any of the following new operators can be used to resolve them:

- *Promotion*: we take an existing instance of the mandatory action $A0$ in the plan and we add the temporal constraint: $S < A0$ if the temporal constraints remain non-contradictory;

- *Adaptation*: add a new instance of the mandatory action *A0* and add the temporal constraint: $S < A0$;
- *Circumvention*: either, add a new step before *S* or add a temporal constraint that would make necessarily true the negation of one of the applicability conditions of the norm;

Dealing with prohibitions. Prohibitions are actions or propositions which are not allowed in a set of intervals, i.e. between any two situations: one where the prohibitions start to be applicable, and another situation where it is no longer prohibited.

*Definition 2 (missing prohibition).* A missing prohibition is a flaw such that the prohibition applicability conditions are necessarily true between the situations *Si* and *Sj*, and:

- The prohibited action is potentially between *Si* and *Sj*;
- or, the prohibited proposition *P* is satisfied in any situation s, such as $Si < s < Sj$.

We introduce the following operators to fix missing prohibitions:

- *Circumvention*: Adding a new step before the *Si*, which would make the negation of one of the applicability conditions necessarily true over the interval [*Si,Sj*];
- *Promotion* or *Demotion*: Add a new temporal constraint which would either (a) move the prohibited action before *Si* (*Promotion*); or (b) move the prohibited action after *Sj* (*Demotion*).

Dealing with permissions. Permissions are actions or propositions which are only allowed under certain conditions. A non-permitted action/proposition is automatically prohibited. To fix missing permissions, we generate for each state a list of prohibitions for all situations or intervals where the permission is not applicable and treat them according to the previous section.

Violating norms. If the agent fails to find a plan through normal planning, we propose that it starts to violate norms starting from the one endowed by the institution he values the least. In MIMOSA, each institution has a priority value that determines its importance. If the planner fails to find a solution when complying with all norms:

- The planner picks randomly a norm in the institution he values the least;
- It returns to the state before the associated normative flaw was resolved and re-plans by ignoring the related normative flaw;
- If the planner still fails, then another additional norm is violated, until a solution is found.

A basic use case of these extensions in renewable resource management would be a villager who ought to provide food for his family, but cannot hunt without a license

in his current area. This raises a flaw, as the villager's plan *to hunt* is in conflict with a prohibition. Using the aforementioned operators, he can either perform a:

- *Promotion* or *Demotion*: he will add steps to get a license before hunting;
- *Circumvention*: he will go elsewhere where he can hunt without a license;
- or else, a *Violation*: he'll consider illegal hunting, i.e. hunting without a license.

### 3.4  The Domain-Specific Language

To build a domain-specific language (DSL) implementing the previous extensions, we use a Model Driven Engineering (MDE) approach with three different artifacts:

(1)  an abstract syntax describing the language's conceptual structure;
(2)  a concrete syntax describing how to write the concepts;
(3)  a semantic describing the meaning of a written sentence.

Abstract Syntax. The abstract syntax describes all the concepts of a DSL and their structure through a metamodel. Since we have already defined how we extend a planning problem with norms through Fig. 2, our DSL's abstract syntax applies these extensions to the abstract syntax of institutions (Fig. 5), and organizations (Fig. 6) by Müller & Raharivelo in [16], with an explicit representation of actions to build agent behavior. An institution is described by [16] as a set of words or concepts, typed by meta-concepts as seen in Fig. 5.

   In summary, an institution is composed of the following elements:

- *meta* (for MetaConcept) enumerates all the meta-concepts;
- *concepts* are categorial concepts that can be denoted to sets of objects through *assertions*, they can also include notions of space or time;



**Fig. 5**  Abstract syntax of an institution according to [16]

**Fig. 6** Abstract syntax of an organization according to [16]

- *indiv* are individual concepts referring to individual objects;
- *norms* describe both regulative and constitutive norms alike.

By using *references*, concepts from other institutions can be imported or inherited. Since our main focus is to generate a sequence of actions, i.e. a plan, we add a STRIPS-like syntax [19] to describe actions by specifying its preconditions and consequences through a set of propositions. Preconditions are propositions that need to be necessarily true to execute the action, while consequences are composed of the propositions describing what the action adds and removes. In summary, the organization's abstract syntax from [16], is represented as in Fig. 6.

Lastly, as we describe the planning problem, we need to describe with the language the abstract syntax of an initial situation and a goal.

Both can be described through a set of propositions which are simple predicates from first-order logic, such as *hasFood()* or *¬hasLicense()*, and can either be an affirmation or a negation (Fig. 7).

Concrete syntax. The concrete syntax depicts how we can use the concepts in the abstract syntax to build valid statements regarding the metamodel, either through some graphics or in some written form.



**Fig. 7** Abstract syntax of a situation and a goal

*Institutions.* The concrete syntax of institutions, added to the definition of a planning problem if Fig. 2 is based upon the concrete syntax of institutions in [16] using the same annotated EBNF (Extended Backus-Naur Form) as follows.

<Institution> ::= 'institution' <institution_name> '{'

   [<references>]
   [<meta>]
   [<concepts>]
   [<indiv>]
   [<norms>]
   [<actions>]

'}'

Unlike the works in [16], our end goal is to produce plans that take norms into account, hence, the actions *component* is added to describe all possible actions which an institution brings about. We explicitly describe actions such as to *cuttingWood* by using a STRIPS representation with preconditions and consequences: the action requires to have an axe, and wood within the current area; executing the action adds the proposition "the agent now has wood", and removes the proposition "there is wood in the current area".

*Organizations.* While institutions can be perceived as virtual entities which define common norms and ontologies, organizations on the other hand are implementations of institutions that denote what will be associated with each concept and can define their own additional concepts if needed.

To describe organizations, we need to specify the institution it is implementing, and which agents or resources play which role through a set of constitutive norms. In terms of concrete syntax, Organizations remain the same as in [16], with the addition of an explicit action representation.

<Organization> ::='organization' <organization_name> from <institution_name>
'{'

   [<meta>]
   [<concepts>]
   [<indiv>]
   [<norms>]

'}'

Situations. Since an initial situation must be provided as a starting point for planning with the language, the following syntax can be used to describe the propositions in each situation.

<Situation> :: = 'situation' <situation_identifier> '{'

   [<Propositions > ].

'}'

Using the initial situation, more situations will be generated internally by the planner as it refines the plan with new actions and constraints to resolve its flaws.

## 4  Discussions

We proposed an approach to extend the POP algorithm with norms through new flaws which will allow agents to modify their plan according to applicable norms. The value of this work is to generate an appropriate agent behavior with regards to norms: either by applying, violating, or circumventing them.

By proposing strategies for an agent to circumvent norms in a spatio-temporal context, we can exhibit new behaviors such as: waiting for the right moment to fish, moving to another space where it is no longer prohibited to have a fishing license, acquiring a social status to be permitted to do something. This expressiveness makes it possible to understand the impact of norms on agent behavior.

The upcoming proof of concept will be done through a replica of the MIRANA model to demonstrate how effective the result is compared to the original ad-hoc implementation. Currently, we still need to take into account the spatial dimension of norms in planning, using spatial algebra like RCC8 (Region Connection Calculus), which has been applied to MAS in [16] to compute all applicable norms given a spatio-temporal context. The next step is to integrate them in our planner by (1) integrating the notions of time and space in situations or intervals of situations; and (2) further refining the proposed normative flaws by specifying which temporal or spatial context causes the normative flaw, and generate behavior accordingly.

While this approach allows us to take norms into account in the behavior generation process, some issues still have to be addressed: (1) the high complexity and performance of the algorithm, (2) the influence of sanctions and rewards on automated planning, and (3) the notion of quantity in planning.

## 5  Conclusion

In this paper, we tackle the issue with agent-based modeling when dealing with norms relative to complex issues such as renewable resource management, where agents are endowed with norms by the organizations in which they play a role in time and space. Since we need a better way to express and account for norms with their temporal, spatial, and social context, the aim of this paper is to propose an agent architecture and its dedicated language's abstract and concrete syntaxes which provides modelers with a tool to model such systems. In that regard, we used automated planning to allow agents to not only take into account or violate norms, but also to consider time and space in their decision-making. Norms are therefore structured in institutions,

and they are instantiated within organizations, which are instances of institutions. The resulting behavior model is built upon the POP paradigm, which allows the agent to exhibit new normative behaviors. The expected architecture, and its language will allow modelers to build one or multiple agents, set a number of norms, and see how they affect an agent's behavior.

## References

1. Aubert, S., Müller, J.-P, Ralihalizara, J.: MIRANA: a socio-ecological model for assessing sustainability of community-based regulations. Int. Environ. Model. Softw. Soc. (iEMSs) 801–808 (2010)
2. Aubert, S., Müller, J.-P.: Incorporating institutions, norms and territories in a generic model to simulate the management of renewable resources. Artif. Intell. Law **21**(1), 47–78 (2013)
3. Müller, J.-P.: The mimosa generic modeling and simulation platform: the case of multi-agent systems. In: Proceedings of the 5th Workshop on Agent-Based Simulation, Lisbon, 2004. SCS. s.l.: s.n., 77–86. International Workshop on Agent-Based Simulation, 5, Lisbonne, Portugal, 3 Mai 2004/5 Mai 2004 (2004)
4. Broersen, J., Dastani, M., Hulstijn, J., van der Torre, L.: Goal generation in the BOID architecture. Cogn. Sci. Q. **2**(3–4), 428–447 (2002)
5. Kollingbaum, M.J., Norman, T.J.: NoA-a normative agent architecture. In: IJCAI, pp. 1465–1466 (2003)
6. dos Santos Neto, B.F., da Silva, V.T., de Lucena, CJ.: Developing goal-oriented normative agents: the NBDI architecture. In : International Conference on Agents and Artificial Intelligence, pp. 176–191. Springer, Berlin, Heidelberg (2011)
7. Rao, A.S., Georgeff, M.P.: Modeling rational agents within a BDI-architecture. In: Fikes, R., Sandewall, E. (eds.) Proceedings of Knowledge Representation and Reasoning (KR&R-91), pp. 473–484. Morgan Kaufmann Publishers, San Mateo, CA (1991)
8. Müller, J.P., Raharivelo, S.O.: Un méta-modèle pour représenter les normes dans un contexte multi-institutionnel territorialisé. Cépaduès (2017)
9. Cialdini, R.B., Trost, M.R.: Social influence: social norms, conformity and compliance. In: D.T. Gilbert, S.T. Fiske, G. Lindzey (eds.) The Handbook of Social Psychology, pp. 151–192. McGraw-Hill (1998)
10. Shoham, Y., & Tennenholtz, M.: On the synthesis of useful social laws for artificial agent societies (preliminary report). In: AAAI, pp. 276–281 (1992)
11. Dignum, V., Vázquez Salceda, J., Dignum, F.O.: Introducing Social Structure, Norms and Ontologies into Agent Organizations. Springer, pp 181–198 (2004)
12. Fornara, N., et al.: Modeling agent institutions. In: Agreement Technologies, pp. 277–307. Springer, Dordrecht (2013)
13. Ostrom, E.: Understanding Institutional Diversity. Princeton University Press (2009)
14. Crawford, S., Ostrom, E.: A grammar of institutions. Am. Polit. Sci. Rev. **89**(3), 582–600 (1995)
15. Interis, M.: On norms: a typology with discussion. Am. J. Econ. Sociol. **70**(2), 424–438 (2011)
16. Raharivelo, S.O., Müller, J.-P.: Un modèle de norme intégrant les conditions spatio-temporelles. In: JFSMA 2018. Systèmes Multi-Agents: Distribution et Décentralisation, pp 117–126. Cé-paduès (2018)
17. Russell, S., Norvig, P.: Artificial Intelligence: A Modern Approach, 2nd edn. Prentice Hall/Pearson Education (2003)
18. Kvarnström, J.: Planning for loosely coupled agents using partial order forward-chaining. In: Twenty-First International Conference on Automated Planning and Scheduling (2011)

19. Fikes, R.E., Nilsson, N.: STRIPS: a new approach to the application of theorem proving to problem solving. Artif. Intell. **2**(3–4), 189–208 (1971)
20. Chapman, D.: Planning for conjunctive goals. Artif. Intell. **32**(3), 333–377 (1987).

# From Threatening Pasts to Hopeful Futures. A Review of Agent-Based Models of Anxiety

**Arvid Horned and Loïs Vanhée**

**Abstract**  Despite being understated, anxiety is a critical factor affecting all levels of society, directly impacting individual decisions and with well-identified ramifications on social play, social constructs, and collective outcomes, as well as being a significant direct social toll tied to yearly trillion-USD social cost. Through a systematic literature review of social simulation research featuring models of anxiety, this paper frames the state of the art on anxiety modelling, and identifies trends and patterns in bibliographic indicators, aspects of anxiety that are modelled, how they are modelled, and their purpose and integration within agent based models. Based on these findings, this paper proposes a way forward as to structure the field as to enable the social simulation community as a whole to cover this critical aspect.

**Keywords**  Anxiety · Social simulation models · Systematic literature review · Psychology models

## 1  Introduction

The trip to the kindergarten has been excruciatingly busy. It is now 08:10, Alice just dropped the kids off and is increasingly anxious about whether she will be in time at her 08:30 meeting. As a result of this anxiety, Alice is likely to become distressed and drive faster; and to arrive exhausted at her meeting, which will affect her interactions with others. Anxiety may also affect the subsequent social play, as Alice's collaborators may offer compassion for Alice's anxiety-driven best (yet risky) efforts to arrive on time, or be annoyed to having had to experience anxiety themselves. Emerging social constructs may also be impacted, for example through a norm: "08:30 meetings are likely delayed", for people to feel less anxious about delayed morning meetings. The ramifications of anxiety can also impact collective outcomes: busi-

A. Horned (✉) · L. Vanhée
Department of Computing Science, Umeå University, SE-901 87 Umeå, Sweden
e-mail: arvid.horned@umu.se

L. Vanhée
e-mail: lois.vanhee@umu.se

nesses might underperform due to ineffective morning activities and eventually shut down. Reciprocally, Alice's anxiety is highly conditioned by personal, environmental, and social constructs: personality, former experiences, culture, values, and norms; and constraints from her business, laws, and road quality. Simulations accounting for anxiety may provide a strong asset for modelling the relation between traffic quality, turnover, and productivity [8], socioeconomic inequalities, and traffic accidents. An anxiety-sensitive simulation may provide a strong argument that the fastest path to a more productive and human-friendly society may be a better road to school.

Anxiety is a psychological process initiated by observing cues of present or future threats triggering response behaviours directed towards dealing with the threat (e.g. preparation, vigilance) [16]. The social impact of anxiety is ubiquitous, spanning across many different domains, including the work-environment, where anxiety impacts health, wellbeing and performance, both positively and negatively [6]; and events such as the global surge of anxiety disorders from 7.3 to 25% during the COVID-19 pandemic [20]. Anxiety reaches US$ 1 trillion yearly cost in social costs from productivity, healthcare, and civil unrest expenses [10]. Anxiety deeply influences social organization as a whole, including status and power allocation, work organization, teaching, and social control, as captured by the uncertainty avoidance cultural dimension [11]. The influence of anxiety on society is further demonstrated by its amplifying effects on the formation of prejudice towards out-groups [4], and trust within groups [13].

However, despite the marked impact of anxiety on society, its inclusion within social simulations is conditioned by our ability to model anxiety accurately and to integrate it successfully within simulations. Besides the challenges and importance of structuring psychology modelling in social simulation research [3, 12], anxiety is inherently particularly difficult to capture (if not, consider) solely relying on intuitions and tacit knowledge: anxiety is a complex and subtle state of mind, which dynamics and ramifications, even when marked, are difficult to adequately seize without extensive psychological expertise and scrutiny. Subsequently, without adequate structuring, the field of simulation is bound to miss opportunities—from disregarding anxiety-sensitive aspects to blank lock-out of high-impact simulation applications; and to produce high-cost simulations that are bound to limited conceptual precision and validation [7].

This paper is dedicated to structuring the current state of the art of anxiety modelling for social simulation. Through a Systematic Literature Review (SLR), this paper scrutinizes the state of the art of Agent-Based Modelling (ABM) and social simulation literature modelling anxiety along four complementary perspectives on modelling: (1) a *bibliographic* overview of the field, assessing the relations between productions, authors' disciplines, venues, and citations over time; (2) *what* components of human anxiety are covered by the various contributions of the field; (3) *how* anxiety is being modelled, covering each step of the design process of social simulation (theory, conceptualization, model, validation, documentation); (4) *why* anxiety is used, relating it to its purpose within the agent deliberation process, social dynam-

Fig. 1 Components of the anxiety framework. Solid lines represent the steps of the anxiety cycle, dashed line represents the influence of factors on other components



Anxiety cycle

ics, and global outcomes based on classic ABM concerns. Then, a more general analysis of the trends of anxiety in social simulation is derived from these results, as well as challenges and opportunities for an effective structuring of the field (Fig. 1).

## 2 Background on Anxiety

Anxiety is tied to multiple experiences including distress, irritation, worry, and sleeping, and concentration issues [1]; to psychological ailments including depression, phobias, panic, eating, and anxiety disorders [9]; to physiological ailments, including cardiovascular, gastrointestinal, and metabolic disorders [19]; to an array of behaviours, including vigilance, and self-sabotage [16]; to altered social plans and constructs such as norms, rituals, and xenophobia [11]. In relation with other psychological constructs, anxiety is commonly differentiated from fear on the basis that anxiety concerns uncertain and anticipated threats, in contrast to fear, which concerns more certain and imminent threats [14]; and from stress in that stress captures a state of sustained arousal that may be caused by many different phenomena including anxiety [1].

Whereas no comprehensive conceptual framework details the process of anxiety formation, the following items appear to be recurrent when crossing the theories provided by the extensive psychological research on anxiety. Anxiety processes can be divided in three steps: **(A) Stimuli**, **(B) Arousal**, **(C) Response**, all being bidirectionally tied to a set of **(D) Factors**. The Stimuli eliciting anxiety is acquired through the **(A1) Perception** of the environment [15] and in particular by the observation of a **(A2) Threat** to one's drives (goals, needs, or values) [15]. As a particular feature of anxiety, stimuli can arise from a **(A3) Future-Oriented** anticipation, i.e., uncertainty or threats regarding plans or other projectives in the future [16]. Second, these stimuli raise the level of (B) arousal, characterized by an **(B1) Intensity**, which is

tied to aversive sensations such as distress, unrest, nervousness, and even panic [15]. This arousal is raised by the perceived degree of **(B2) Uncertainty** concerning the threat, including the likelihood and consequences of the threat and believed personal vulnerability [5]. Moreover, anxious arousal can trigger **(B3) Sensitization** to further anxiety-inducing stimuli in the environment, increasing the importance associated to further associated anxiety-inducing stimuli (confirming the threat) [15]. A (C) Response to anxiety usually involves either dealing with the threat, usually through increased control, or reducing the feeling of anxiety, through coping. Seeking **(C1) Pragmatic control** over the threat, by taking concrete steps to confront, mitigate, and avoid the threat and its effects [16]. Seeking **(C2) Epistemic control** directed towards reducing uncertainty regarding the occurrence and the consequences of the threat, with information seeking behaviour [16]. Last, **(3) Coping** seeks to reduce the negative experiences tied to anxiety, such as distracting oneself or displacing the anxiety, which is a common source of maladaptive responses to anxiety as it can override functional action [9]. Last, anxiety is tied to external (D) Factors, including individual **(D1) traits** that influence one's sensitivity to various forms of anxiety, such as genetic dispositions to experience anxiety or situational dispositions (e.g. age, social status) [1]. Anxiety is particularly sensitive to **(D2) Learning** as one's assumed ability to sustain control and past traumas influence other steps of the anxiety process [15]; and **( D3) Long-term consequences** of experiencing anxiety, mostly tied to the development of anxiety-disorders [9].

## 3   Method

To answer our research question, relevant texts were identified, screened, reviewed and analysed following a Systematic Literature Review (SLR) process utilizing the Prisma checklist [18]. The phases of the SLR and included papers are described in Fig. 2. First, papers were collected using the search strategy and query described in Table 1. To identify relevant literature, the collected papers were screened by the first author to include agent-based models in which anxiety is integrated. After the screening process, the full texts of the remaining papers were subjected to a systematic content analysis following the codes and patterns described in the codebook. Cited agent-based models of anxiety not captured by the search strategy were added to the SLR during the coding phase. The codebook was developed using a combination of inductive and deductive coding, first establishing a ground set of codes to identify patterns, then adding new codes as new patterns emerged during the analysis.

First, to produce a bibliographical overview, the evolution of the field over time was investigated by coding the meta-data, with codes including: (1.a) Publication Year; (1.b) Publication outlet; (1.c) Publication outlet domain. Moreover, the affiliation of the authors and citations was investigated to the build up of models within scientific disciplines with codes including: (1.d) What are the authors' scientific disciplines; (1.e) What earlier models of anxiety do they cite?

**Fig. 2** The different phases of the systematic literature review. Search date: 11.05.2022

**Table 1** Search strategy and inclusion criteria in the systematic literature review

| | |
|---|---|
| Search strategy | Scopus database (www.scopus.com) was selected as the primary source, after comparison |
| Search query | TITLE-ABS-KEY (("social simulation" OR "agent based model") AND "anxiety") |
| Inclusion criteria | Paper is accessible (authors were contacted if no reachable link)<br>Paper is in English<br>Paper presents an agent-based model or a social simulation<br>Anxiety is modelled by the model or simulation |

Second, a review of what aspects of anxiety were included into the model of anxiety was assessed based on the framework provided in Fig. 1, coding the level of implementation using four labels: (None) There is no connection to this component; (Interpretable) The inclusion of this component can be interpreted to exist in the model; (Implicit) This component is included implicitly with concepts directly connected to it; (Explicit) This component is explicitly included into the model. The implementation of anxiety components were coded as follows: (2.a) A1: Perception; (2.b) A2: Threat; (2.c) A3: Future orientation; (2.d) B1: Intensity of arousal caused by anxiety; (2.e) B2: Uncertainty; (2.f) B3: Sensitization; (2.g) C1: Pragmatic Control; (2.h) C2: Epistemic Control; (2.i) C3: Coping; (2.j) D1: Trait; (2.k) D2: Learning;

(2.l) D3: Long-term consequences. As an example, threat (A2) is coded as explicit if the model of anxiety includes a trigger of anxiety that represents threats to the agent, and uncertainty (B2) is coded as implicit if the model of anxiety concerns an uncertainty but does not include a representation of it.

Third, to assess how former papers model anxiety, their process was coded following the steps of model development [17], including: (3.a) What theoretical foundation (and whether it is grounded in psychology) is used to model anxiety?; (3.b) Is the model of anxiety conceptualized?; (3.c) Is the model based on any experiment or data?; (3.d) How is anxiety represented in the model?; (3.e) How is individual variability in anxiety represented? (3.f) How is the model of anxiety validated?; (3.g) How is the model of anxiety documented?

Lastly, why anxiety is modelled was investigated by assessing the role anxiety plays in the larger context of the paper, with questions including: (4.a) What is the model application domain?; (4.b) How is anxiety integrated?; (4.c) What is the motivation for including anxiety?; (4.d) Is anxiety used as a proxy for a separate concept? To relate anxiety with the types of deliberations made by the agent, we checked whether they include the psychological features identified by Jager's EROS [12], and whether it is directly connected to anxiety with codes including: (4.e) A.I: Theory of normative conduct; (4.f) A.II: Goal Frame theory; (4.g) A.III: Similarity theory; (4.h) A.IV: Social judgement and opinion dynamics; (4.i) A.V: Elaboration likelihood model; (4.j) A.VI: Theory of planned behaviour; (4.k) A.VII: Integrated models. In the same fashion, the social features in agent deliberation were investigated by applying the social features identified by Bourgais [2], coding: (4.l) B.I: Cognition; (4.m) B.II: Personality; (4.n) B.III: Emotions; (4.o) B.IV: Social relations.

## 4  Results

Of the 32 papers retrieved by the query on the Scopus database, 20 papers were excluded due to: no access to paper (n = 1); no ABM or social simulation was presented (n = 10), anxiety was not included into the model (n = 9). The citation analysis resulted in two additional papers being added, resulting in a total of 14 papers included in the final analysis (found in Annex).

The **_Bibliographical overview_**, included in Fig. 3, depicts the number of papers published per year, this figure shows that modelling anxiety within ABMs is a relatively new phenomenon with the first retrieved paper dating from 2009. The scientific productivity over the field is relatively limited but steady, ranging from one to three papers yearly. Authorship is relatively interdisciplinary, 57% (n = 8) of the papers are written by authors from different disciplines. More than half of the papers presented their model of anxiety (n = 9) without any reference to any previous model of anxiety, with a small number (n = 5) of papers building on models that had been published previously.

**Fig. 3** Number of publications per year, with colours representing the application domain of the publications

*Components of anxiety* identified in Sect. 2, and the degree to which they are modelled in the included literature, is displayed in Fig. 4.

**Stimuli** triggering anxiety, including the perception (A1) of it, is modelled explicitly in more than half of the papers (n = 9). Threats (A2) are modelled explicitly in roughly 75% of the papers (n = 11), with anxiety-increasing inputs representing aversive events to the agent. However, no model includes future-orientated (A3) components of anxiety (n = 0), with all sources and behaviours associated with anxiety only existing in the present.

**Arousal** and the intensity (B1) of it is modelled explicitly in roughly 64% of the papers (n = 9). Uncertainty (B2) is modelled implicitly in just over half of the papers (n = 8), and in less than half of the papers it can be interpreted (n = 5), modelling sources of anxiety that may be connected to uncertainty. No models include sensitization (B3) (n = 0), with current anxiety-levels having no influence on the anxiety generated in response to later threats.

**Responses** to anxiety, such as pragmatic control (C1) is modelled explicitly in just under half of the papers (n = 6), with agents taking anxiety-reducing actions that target the problem at the source of the anxiety. In contrast, no models include behaviour to seek epistemic control (C2) (n = 0), such as information-seeking behaviour in response to anxiety. Coping (C3) behaviour is modelled explicitly in roughly a quarter of the papers (n = 4), with agents engaging in behaviour to alleviate the experience of anxiety without affecting the threat.

**Factors**, including individual traits (D1) affecting the degree of arousal raised by anxiety, could be interpreted in less than a quarter of the models (n = 2), as the models include some inherent heterogeneity between agents influencing anxiety formation. No models include learning (D2) (n = 0) in relation to sources of anxiety or behaviour in response to it. Moreover, no models integrated any long-term consequences (D3) of anxiety.

What components of anxiety are modelled?

| ID | Application Domain | Anxiety Representation | A1 | A2 | A3 | B1 | B2 | B3 | C1 | C2 | C3 | D1 | D2 | D3 |
|----|--------------------|----------------------|----|----|----|----|----|----|----|----|----|----|----|----|
| 1 | Primate social behavior | Scale-based | | | | | | | | | | | | |
| 2 | Primate social behavior | Scale-based | | | | | | | | | | | | |
| 3 | Primate social behavior | Scale-based | | | | | | | | | | | | |
| 4 | Primate social behavior | Scale-based | | | | | | | | | | | | |
| 5 | Primate social behavior | Scale-based | | | | | | | | | | | | |
| 6 | Electric vehicle driving and charging | Constant-based | | | | | | | | | | | | |
| 7 | Virtual Interview Coach | Undefined | | | | | | | | | | | | |
| 8 | Building evacuation | Undefined | | | | | | | | | | | | |
| 9 | Religiosity | Scale-based | | | | | | | | | | | | |
| 10 | Religiosity | Scale-based | | | | | | | | | | | | |
| 11 | Electric vehicle driving and charging | Scale-based | | | | | | | | | | | | |
| 12 | Reaction to climate information | Scale-based | | | | | | | | | | | | |
| 13 | Reaction to climate information | Scale-based | | | | | | | | | | | | |
| 14 | Electric vehicle driving and charging | Constant-based | | | | | | | | | | | | |

**Fig. 4** Components of anxiety from the framework in Fig. 1, darker shades represent higher levels of inclusion into the model of anxiety. Crossed with the application domain of the model and how anxiety was represented

*How anxiety is modelled* is displayed in Fig. 4. Two clusters emerge out of the data when looking at what aspects of anxiety are being modelled.

**Theoretical foundations** are grounded in a psychological theory in 36% of the papers (n = 5). Terror Management Theory (TMT) is used by two papers as a foundation for modelling anxiety, a theory describing how anxiety about death may facilitate engagement in religious belief and rituals. Two other papers use the Threat and Defence Theory as a basis, describing how individuals may either approach or distance themselves from a threat based on individual motivations. Lastly, one paper uses Flow Theory as a foundation, where anxiety is described as a state of low motivation triggered by a perceived lack of competence. Most frequently, models of anxiety are based on general theories developed in the field of the simulated context (n = 9) including: range anxiety concerning the uncertain distance one can travel by car (n = 3), anxiety in evacuation of buildings during earthquakes (n = 1), anxiety displayed by primates (n = 5).

**Conceptualization** of anxiety is not to be found. While papers may provide a conceptual framework of the modelled phenomena which includes anxiety, no papers provide a conceptual framework of the anxiety they include into their models (n = 0). All papers directly include anxiety as a component of the model, without arguing for the dynamics. However, all models referred to earlier experiments or data reporting on anxious phenomenon (n = 14).

**Modelling** anxiety, 12 papers introduced formal models of anxiety, other papers only informally explained its use. Models of anxiety can be clustered in two categories: *constant-based models* (n = 2) and *scale-based models* (n = 10). In constant-based models, anxiety is represented as a constant tied to the agent that plays as a modifier in decisions. In scale-based models, anxiety changes over time, being recorded by a [0, 1] variable, representing the intensity of experienced anxiety. This variable is updated every round, increased when the agent faces uncertainty or threats and decreased if the agent performs a certain action. In terms of influence, scale-based models tie higher levels of anxiety to triggering dedicated anxiety-reducing actions, either directly (e.g. applying a specific pre-set anxiety-coping action if the trigger

is met and applicable) or indirectly (e.g. increasing the odds of triggering an action effectively lowering anxiety).

**Validation and documentation**: No paper included a dedicated validation of the anxiety model. One paper documented how anxiety was modelled, focusing on the update process in response to perceived hazards in the environment.

***Purpose and Integration of Anxiety in ABMs*** is covered by the simulation domain, and how anxiety is integrated and tied with agent decision aspects.

**Simulation domain**. Six application domains are covered by the literature including: electric vehicle adoption (n = 3), religiosity (n = 2), earthquake building evacuation (n = 1), reaction to negative information about climate change (n = 2), social behaviour in primates (n = 5), and motivation in coaching (n = 1).

**Anxiety integration**, in roughly two thirds of the papers, is a driving force behind a particular behaviour in the simulated context (n = 10), and also motivated the inclusion of anxiety as such (n = 10). In two other cases, anxiety was integrated as a global variable affecting the decision-making of agents in the simulated context (n = 2), and once as a mental state emerging from other processes in the agent (n = 1). Modelling anxiety was motivated in almost all the papers (n = 12) where it was not a proxy for something else (n = 12). In the remaining two cases, anxiety is used as a proxy for a low motivational state (n = 1) and a state of panic (n = 1). The homeostasis of anxiety is integrated as a decline whenever the source of the anxiety was not present (n = 7) in half of the models, in two models actions are the only way to reduce anxiety (n = 2), and in one case anxiety completely dissipates as the source of anxiety disappears (n = 1).

**Agent deliberation aspects** include the psychological features described in [12], anxiety played a frequent role in the goal framing of agents (n = 10), influencing the agents' motivation to engage in a particular behaviour. Modelling planned behaviour including intention to perform a particular behaviour influenced by personal attitudes, norms and behavioural control was less common (n = 6), but was often influenced by anxiety (n = 5). Normative conduct was included in 28.6% (n = 4) of models. However, anxiety only influenced normative conduct in one model, where anxiety increased the likelihood of influencing other agents forming norms. Two models included similarity theory (n = 2) and social judgement (n = 2) as part of agent interaction, but did not include the influence of anxiety on these aspects. No model made a distinction between central and peripheral processing, building on the elaboration likelihood model. Similarly, no model presented an integrated model of these components. Cognition and reasoning with mental states was included in 35.7% (n = 5) of models, with anxiety influencing cognition in three of the models as it drove confirmation of held beliefs (n = 2) or arose from held beliefs (n = 1). No models incorporated personality as an influence on deliberation. Three models incorporated a range of emotions influencing the deliberation of agents, where anxiety was represented as one of the emotions. Half of the models incorporate social relationships influenced by anxiety (n = 7) as it drove behaviour either strengthening or harming the affiliation between agents (Fig. 5).

**Psychological and social features of the ABM**

| ID | Anxiety Representation | Anxiety integration | A.I | A.II | A.III | A.IV | A.V | A.VI | A.VII | B.I | B.II | B.III | B.IV |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Scale-based | Anxiety driving force | | | | | | | | | | | |
| 2 | Scale-based | Anxiety driving force | | | | | | | | | | | |
| 3 | Scale-based | Anxiety driving force | | | | | | | | | | | |
| 4 | Scale-based | Anxiety driving force | | | | | | | | | | | |
| 5 | Scale-based | Anxiety driving force | | | | | | | | | | | |
| 6 | Constant-based | Global parameter | | | | | | | | | | | |
| 7 | Undefined | Emergent state | | | | | | | | | | | |
| 8 | Undefined | Undefined | | | | | | | | | | | |
| 9 | Scale-based | Anxiety driving force | | | | | | | | | | | |
| 10 | Scale-based | Anxiety driving force | | | | | | | | | | | |
| 11 | Scale-based | Anxiety driving force | | | | | | | | | | | |
| 12 | Scale-based | Anxiety driving force | | | | | | | | | | | |
| 13 | Scale-based | Anxiety driving force | | | | | | | | | | | |
| 14 | Constant-based | Global parameter | | | | | | | | | | | |

**Fig. 5** The inclusion of psychological and social features of agent deliberation, with three shades representing: (light) feature is not included into agent deliberation; (medium) feature is included into agent deliberation; (dark) feature is included into agent deliberation with a connection to anxiety. Features include: (A.I) Theory of normative conduct; (A.II) Goal frame theory; (A.III) Similarity theory; (A.IV) Social judgement (A.V) Elaboration likelihood model; (A.VI) Theory of planned behaviour; (B.I) Cognition; (B.II) Personality; (B.III) Emotions; (B.IV) Social relations

## 5 Discussions, Conclusions and Way Forward

Summing up the results, anxiety in social simulation fits the profile of an emerging topic that fosters simultaneous interest across disconnected communities, leading to the re-invention of near-identical models with similar advantages and shortcomings. Overall, anxiety in social simulation is a relatively young topic, receiving limited but steady interest scattered across disciplines. Despite simultaneous inventions, models being built have overall very similar features, (1) in the aspects of anxiety they cover (stimuli, threats, intensity, uncertainty, pragmatic control, coping) and fail to cover (future-orientation, sensitization, epistemic control, learning, long-term consequences of anxiety); (2) in the modelling process (occasional grounding in psychological theories, no conceptualization, occasional explicit description of the model, no validation, no documentation); (3) in the technicalities of the model (two types of models: constant-based and scale-based); (4) in the aspects of agent deliberations anxiety is connected with (goals and social interactions) and non-connected with (similarity theory, social judgement, personality, emotions), and (5) with similar function in terms of agent behaviour (as a factor influencing a single decision, for triggering a specific behaviour following repeated exposure to specific threats). A contrario, recovered items also show a wide variety of domains (e.g. electric vehicle, religion, primate behaviour)—highlighting the acknowledgement of the breadth of applicability of anxiety. Two main trends of approaches emerge, correlating methods, theories, models, and outcomes: a *fine-grained approach*: theory-based, scale-based models covering threat sensitivity, perception and either pragmatic control or coping behaviours, and tied to goals, and either plans or social relations; and a *coarse approach*, reducing anxiety as a modifier for a decision tied to uncertainty, thus representing trait anxiety.

This overview shows that now is the time to structure the field of anxiety in social simulation, as despite its initial success, the intuitionist ad-hoc approach is reaching its limits. First, the set of models generated is shown to reach hard limits: while having their virtues, these models are also blind to some primordial shapes of anxiety (future orientation, social judgement, emotions, anxiety bleeding over multiple decisions) that can be modelled [21] and their ramifications on behaviour, interactions, and collective outcomes—in addition to inherent grounding and validation limits of the theories, concepts, and models. Second, the ad-hoc approach is also limiting in terms of *methods and applications*, as this approach fails to study the alternative theories, conceptualizations, and models; nor develop a pragmatic general assessment of the relevance of the inclusion of anxiety in the model. For both reasons, further structuring appears necessary for escaping stagnation, as to access the landscape of high-impact simulations of situations in which anxiety is a significant factor. As to alleviate these two limitations, the following two challenges bring forward two complementary courses of action as to foster a structured, incremental development of the field.

**Challenge 1** *Develop dedicated models of anxiety, including theoretical reviews, justified conceptualizations, and validated models.*

As to move beyond the status quo, we need to develop incremental theories, conceptualizations, models, and validations through dedicated research, beyond the current approach in which anxiety is an integrated component that cannot be discussed, justified and validated, following the general approach for producing high-standards models arising from psychology [7, 12]. Concrete actions include: (1) structuring theories of anxiety in light of the aspects of anxiety they cover; (2) developing detailed conceptual frameworks of anxiety (e.g. ontologies, processes); (3) developing generic models of anxiety; (4) develop implementations of these models and make them available to the community; (5) develop validation methodologies and validate these models of anxiety. In addition, we invite readers to contribute to a database of models of anxiety, taking a direct collective action towards structuring the field[1]

**Challenge 2** *Develop impact assessment frameworks and methods for applying anxiety-sensitive modelling in practice.*

How much, and causal chains explaining how anxiety impacts society are difficult to identify solely based on intuitions, therefore solutions for pragmatic assessment and understanding of the dynamics of anxiety is required for the accurate integration of anxiety within simulations. Concrete actions include: (1) identifying high-impact or clear-cut applications of social simulations, through identifying situations and social constructs strongly tied to anxiety; (2) identifying factors coupling situations to anxiety, developing assessment tools; (3) develop causal chains relating anxiety to behaviours, interactions, and collective outcomes; (4) identifying relations between classic ABM components and anxiety (e.g. trust, norm formation); (5) developing and documenting simulations of anxiety-sensitive applications.

---

[1] We invite readers to share social simulation models which anxiety is a component of. You can find the database along with instructions using the following link: http://s.cs.umu.se/i1d2bq.

Last, both challenges are to be crossed: developing methodologies for including anxiety within models, connecting anxiety factors raised by the situation to technical solutions for modelling these factors. By blending a clear understanding of anxiety and its ramifications on society and developing advanced and reliable models of anxiety and its integration within agent deliberation, we will develop the incremental ability for capturing high-stakes social situations with our simulations, as well as having a new tool for revisiting and improving former approaches. As a playground for forecasting the future of societies and acting upon them, social simulations that accurately replicate human-like anxiety provide us with unprecedented insights into a strong psychological driving force that is critical to wellbeing, giving us an edge for creating more hopeful futures for human societies. In an age of major change, sustaining hope about the future may become key insights that social simulation can offer.

## Annex

ID:1  I. Puga-Gonzalez, H. Hildenbrandt, and C. K. Hemelrijk, 'Emergent Patterns of Social Affiliation in Primates, a Model'

ID:2  C. K. Hemelrijk and I. Puga-Gonzalez, 'An individual-oriented model on the emergence of support in fights, its reciprocation and exchange'

ID:3  E. Evers, H. De Vries, B. M. Spruijt, and E. H. M. Sterck, 'The EMO-model: An agent-based model of primate social behavior regulated by two emotional dimensions, anxiety-FEAR and satisfaction-LIKE'

ID:4  E. Evers, H. De Vries, B. M. Spruijt, and E. H. M. Sterck, 'Emotional bookkeeping and high partner selectivity are necessary for the emergence of partner-specific reciprocal affiliation in an agent-based model of primate groups'

ID:5  E. Evers, H. Vries, B. M. Spruijt, and E. H. M. Sterck, 'Intermediate-term emotional bookkeeping is necessary for long-term reciprocal grooming partner preferences in an agent-based model of macaque groups'

ID:6  C. J. R. Sheppard, A. R. Gopal, A. Harris, and A. Jacobson, 'Cost-effective electric vehicle charging infrastructure siting for Delhi'

ID:7  N. S. Ajoge, A. A. Aziz, and S. A. Mohd Yusof, 'On modeling of interviewee motivation mental states for an intelligent coaching agent'

ID:8  G. P. Cimellaro, F. Ozzello, A. Vallero, S. Mahin, and B. Shao, 'Simulating earthquake evacuation using human behavior models'

ID:9  F. L. Shults, J. E. Lane, W. J. Wildman, S. Diallo, C. J. Lynch, and R. Gore, 'Modelling terror management theory: computer simulations of the impact of mortality salience on religiosity'

ID:10  F. L. Shults, R. Gore, W. J. Wildman, C. J. Lynch, J. E. Lane, and M. D. Toft, 'A generative model of the mutual escalation of anxiety between religious groups'

ID:11 M. Van Der Kam, A. Peters, W. Van Sark, and F. Alkemade, 'Agent-based modelling of charging behaviour of electric vehicle drivers'

ID:12 M. L. Kapeller and G. Jäger, 'Threat and anxiety in the climate debate: an agent-based model to investigate climate scepticism and pro-environmental behaviour'

ID:13 M. L. Kapeller, G. Jäger, and M. Füllsack, 'Social Norms and the Threat of Climate Change: An Agent-Based Model to Investigate Pro-Environmental Behaviour'

ID:14 T. Gnann, D. Speth, K. Seddig, M. Stich, W. Schade, and J. J. Gómez Vilchez, 'How to integrate real-world user behavior into models of the market diffusion of alternative fuels in passenger cars—An in-depth comparison of three models for Germany'

# References

1. Barlow, D.H.: Anxiety and Its Disorders. Guilford press (2004)
2. Bourgais, M., Taillandier, P., Vercouter, L.: BEN: an architecture for the behavior of social agents. J. Artific. Soc. Soc. Simul. **23**(4), 12 (2020). https://doi.org/10.18564/jasss.4437
3. Bourgais, M., Taillandier, P., Vercouter, L., Adam, C.: Emotion modeling in social simulation: a survey. J. Artific. Soc. Soc. Simul. **21**(2) (2018). https://doi.org/10.18564/JASSS.3681
4. Brader, T., Valentino, N.A., Suhay, E.: What triggers public opposition to immigration? anxiety, group cues, and immigration threat. Am. J. Pol. Sci. **52**(4), 959–978 (2008). https://doi.org/10.1111/j.1540-5907.2008.00353.x
5. Carleton, R.N.: Into the unknown: a review and synthesis of contemporary models involving uncertainty. J. Anxiety Disord. **39**, 30–43 (2016). https://doi.org/10.1016/j.janxdis.2016.02.007
6. Cheng, B.H.: Understanding the dark and bright sides of anxiety: a theory of workplace anxiety. J. Appl. Psychol. **103**(5), 537 (2018). https://doi.org/10.1037/APL0000266
7. Gabriel, G.T., Campos, A.T., Leal, F., Montevechi, J.A.B.: Good practices and deficiencies in conceptual modelling: a systematic literature review. J. Simul. **16**(1), 84–100 (2022). https://doi.org/10.1080/17477778.2020.1764875
8. Gobind, J.: Transport anxiety and work performance. SA J. Human Res. Manag. **16**(1), 1–7 (2018). https://doi.org/10.4102/sajhrm.v16i0.943
9. Grupe, D.W., Nitschke, J.B.: Uncertainty and anticipation in anxiety: an integrated neurobiological and psychological perspective. Nat. Rev. Neurosci. **14**(7), 488–501 (2013). https://doi.org/10.1038/nrn3524
10. Health, T.L.G.: Mental health matters. Lancet Global Health **8**(11), e1352 (2020). https://doi.org/10.1016/S2214-109X(20)30432-0
11. Hofstede, G.: Dimensionalizing cultures: the Hofstede model in context. Online Read. Psychol. Culture **2**(1) (2011). https://doi.org/10.9707/2307-0919.1014
12. Jager, W.: Enhancing the realism of simulation (EROS): on implementing and developing psychological theory in social simulation. J. Artific. Soc. Soc. Simul. **20**(3), 14 (2017)
13. Kenworthy, J., Jones, J.: The roles of group importance and anxiety in predicting depersonalized ingroup trust. Group Proc. Intergroup Relat. **12**(2), 227–239 (2009). https://doi.org/10.1177/1368430208101058

14. Lang, P.J., Davis, M., "Ohman, A.: Fear and anxiety: animal models and human cognitive psychophysiology. J. Affect. Disord. **61**, 137–159. (2000). https://doi.org/10.1016/S0165-0327(00)00343-8

15. Levy, I., Schiller, D.: Neural computations of threat. Trends Cognit. Sci. **25**(2), 151–171 (2021). https://doi.org/10.1016/j.tics.2020.11.007

16. Miceli, M., Castelfranchi, C.: Anxiety as an "epistemic" emotion: an uncertainty theory of anxiety. Anxiety, Stress Coping **18**(4), 291–319 (2005). https://doi.org/10.1080/10615800500209324

17. Nikolic, I., Ghorbani, A.: A method for developing agent-based models of socio-technical systems. In: 2011 International Conference on Networking, Sensing and Control, pp. 44–49 (2011). https://doi.org/10.1109/ICNSC.2011.5874914

18. Page, M.J., McKenzie, J.E., Bossuyt, P.M., Boutron, I.H., et al: The PRISMA 2020 statement: an updated guideline for reporting systematic reviews. PLOS Med. **18**(3), e1003,583 (2021). https://doi.org/10.1371/journal.pmed.1003583

19. Sala, T., Cox, B.J., Sareen, J.: Anxiety disorders and physical illness comorbidity: an overview. In: Anxiety In Health Behaviors and Physical Illness, pp. 131–154. Springer New York (2007)

20. Santabárbara, J., Lasheras, I., Lipnicki, D.M., Bueno-Notivol, J., Pérez-Moreno, M., López-Antón, R., De la Cámara, C., Lobo, A., Gracia-García, P.: Prevalence of anxiety in the COVID-19 pandemic: an updated meta-analysis of community-based studies. Prog. Neuro-Psychopharmacol. Biol. Psychiatry **109**, 110–207 (2021). https://doi.org/10.1016/j.pnpbp.2020.110207

21. Vanhée, L., Jeanpierre, L., Mouaddib, A.I.: Anxiety-sensitive planning: from formal foundations to algorithms and applications. In: International Conference on Automated Planning and Scheduling. Singapore (2022)

# Impact of Leader-Follower Behavior on Evacuation Performance: An Exploratory Modeling Approach

**Jakob Irnich, Natalie van der Wal, Dorine Duives, and Willem Auping**

**Abstract** Different leader-follower behaviors may be observed in models, such as group gathering, backtracking, and changing between groups. However, a comparison of these behaviors resulting in possible substantially different estimates of optimal evacuation procedures is lacking. Hence, we developed an agent-based model in combination with exploratory modeling to compare backtracking, group gathering, and followers changing leaders and investigate their influence on the evacuation and response time. The simulation results showed that backtracking and changing of groups increased the evacuation time. Whereby group gathering increase the response time. In addition, the combination of behaviors increases the influence on evacuation and response time. Further research needs to test these results with empirical studies and investigate the impact of other leader-follower behavior. The found insights may be utilized in evacuation research for modeling this behavior and they provide a valuable basis for designing policies in buildings with a high distribution of leader-follower groups.

**Keywords** Leader-following behaviour · Evacuation · Agent-based modelling · Uncertainty · EMA workbench

## 1 Introduction

In general, leadership can be seen as a core attribute of social groups [1]. Haghani et al. found that leadership was one of the most influencing decision-making processes during an evacuation [2]. In addition, real-life observations of evacuations revealed

J. Irnich (✉) · N. van Wal · W. Auping
Faculty of Technology Policy and Management, Delft University of Technology,
Jaffalaan 5, 2628 BX Delft, The Netherlands
e-mail: J.V.Irnich@student.tudelft.nl

D. Duives
Faculty of Civil Engineering and Geosciences, Delft University of Technology,
Stevinweg 1, 2628 CN Delft, Netherlands

that leaders play an essential role in the evacuation process [3, 4]. The leader drives the group's movement and thus influences the follower in his decisions towards the exit [5]. Various researchers have already explored leader-follower behavior with the help of empirical studies [6, 7]. For instance, Jones and Hewitt [3] realized that a leader might be imposed through hierarchical structures or emerge spontaneously. In addition, the group may split in case of different opinions, resulting in a new group with another leader. In line with observations in empirical studies, researchers implemented leader-follower behavior in models. For instance, Li et al. [8] developed a social force model, including leader-follower behavior. Other authors incorporated that the group is gathering before the evacuation [9]. Yet, leader-follower behavior implemented in evacuation models differ substantially. These differences in model implementation potentially result in different estimates of the optimal evacuation procedure. A thorough comparison of different model implementations is essential to better understand the impact of leader-follower behaviour models on the evacuation performance. This research aims to determine how three different leader-follower behaviours influence the evacuation and response time in buildings, namely back-tracking, group gathering and followers changing leaders.

The remainder of this paper first presents the methodology in Chap. 2. Chapter 3 introduces an innovative Agent-based model and provides verification and validation. After the model presentation, the results are shown in Chap. 4. Finally, the article ends with a discussion of the results and conclusion in Chap. 5.

## 2 Methodology

In order to identify the effects of the three leader-follower behaviors, we first need to develop a suitable model and then establish experiments to receive a robust result for the influence. We used 2 distinct methodologies, namely Agent-based (ABM) and exploratory modeling. Below, both methods are briefly outlined.

### 2.1 Agent-Based Modeling

Various methodologies exist to model evacuations, such as social force models, fluid dynamics, and ABM [10]. Each methodology may be utilized for unique research goals. As the research investigates different behaviors and their influence on the emergent pattern in a complex environment, ABM is a suitable methodology for this study. Due to its bottom-up approach and ability to incorporate flexible and autonomous actions of agents in an environment [11], ABMs enable the integration of evacuee relationships and building interactions during an evacuation. Especially these attributes lead to choosing an ABM.

## 2.2 Exploratory Modeling

Ronchi et al. [12] identified four different uncertainties, which are predominant in evacuation research: Input, measurement, behavioral and intrinsic uncertainty. In an exploratory analysis, different parameter combinations in the parameter space will be chosen in order to investigate how the model behaves under the influence of uncertainties. Exploratory models do not predict or find precise answers to specific questions [13]. However, it develops insights regarding the behaviour of the model and helps discover extreme model behaviors [13]. For instance, feature scoring, explores the influence of uncertainties of the model. Higher confidence in results and thus a more robust solution may be achieved [14]. For leader-follower behaviors in evacuations, exploratory modeling may accomplish robust results about the influence on evacuations, independent of one particular scenario, increasing the overall value of this study for the evacuation research community.

## 2.3 Simulation Procedure

The difference between traditional and exploratory modeling is the absence of a base case but the utilization of a base ensemble [15]. A base ensemble represents a sample over the uncertainty space. In our model, we used a Latin hypercube sampling with 1000 scenarios. In addition, feature scoring may help identify the relevance of uncertainties on the KPIs [16] and is thus applied. Finally, we performed a multivariant behavior testing on the base ensemble. The key performance indicators of interest are total evacuation, and the mean response time. Whereby the total evacuation time is defined by the time between the start of the evacuation and the last agent leaving the building and the response time is determined by the time period between the recognition time of the agent and the first movement towards the exit. Furthermore, to study the behavior inside the groups and how this varies with additional policies, we monitored as a secondary outcome the mean intragroup distance between the groups. We conducted all experiments with the Exploratory Modelling and Analysis (EMA) Workbench. A detailed description of the EMA Workbench may be found in [14].

# 3 Model Representation

We developed an ABM, including different leader-follower behavior. A detailed description of the model is found on request at https://github.com/JIRnic. First, the purpose of the model is explained, then the state variables and states are shown [17]. The next part describes the process, the leader-follower behavior and the uncertainties. Finally, the chapter finishes with a short explanation of verification and validation.

## 3.1   Purpose of the Model

The purpose of the model is to investigate how specific leader-follower behavior may influence the evacuation and response time inside buildings. In particular, we examined three behaviors in more detail: backtracking, group gathering before the evacuation, and followers to change to another leader. Our goal is to receive a robust result regarding these three "policies" in buildings with the help of exploring the uncertainty space in an evacuation process.

## 3.2   State Variables and States

Overall, our agent-based model consists of three hierarchical levels: the individual level of each agent, the spatial level and the environment [18].

**Agents**: Our model contains three different agents: the leaders, followers, and individual evacuees. Leaders and followers are members of a group. Whereby the leader searches the path, and the follower follows the leader. Individual agents evacuate on their own.

**Layout**: We designed the layout of the model to mimic a museum or municipality hall. It contains five exits, whereby the main exits are located on the left and right sides of the main hall. The three other emergency exits are positioned at the top and bottom. The width of each exit is set to two meters. Black cells represent walls and obstacles, which must be avoided by agents. The building is illustrated in Fig. 1. A symmetrical layout was utilized to minimize the influence on the evacuation performance of where groups and individuals are placed.

**Operationalization layout**: The environment includes the scale and time dimension of the model. The software Netlogo represents the environment as a grid, in which one patch represents an area of $1 \times 1$ m in real-life. In addition, time is epitomized by ticks. For each tick, an agent is following specific rules. In the model, one tick symbolizes one second in the real time.

## 3.3   Process Overview and Scheduling

We divided the model into three phases, namely the pre-movement, movement and queuing [19]. The pre-movement step may further be subdivided into the recognition and response time [19]. After the pre-movement process, agents move towards the exit. Before leaving the building, the agent needs to queue as the door may be blocked by other agents. The three different Leader-follower behavior may be additionally added to the pre-movement and movement phase. The overall high-level process is shown in Fig. 2.

**Fig. 1** Building layout



**Fig. 2** The overall process for each agent in the model

## 3.4  Leader-Follower Behavior

We included and explored three different kinds of leader-follower behaviors: Back-tracking, Group gathering before the evacuation, as well as Flexibility of the group.

**Backtracking**: Backtracking is a behavior that has been observed by leaders in social groups, for instance, close friends and family [20]. The group members evacuate and try to stay together throughout the whole evacuation process [20]. However, a group member may depart from the group in the rush of the evacuation and interaction with other evacuees [21]. In order to reestablish the connection, the leader reduces its speed and delays its evacuation until the lost member has caught up [21].

**Group gathering**: Group members may perform different actions during the recognition and response phase. After every group member finishes their task, social groups gather before evacuating together [3]. In our model, every group member moves towards the leader. Only if all group members are within a range defined by a threshold the leader starts evacuating. Here, we operationalized the threshold using the work of Moussaïd et al. [5]. As Group gathering may be allocated to the pre-movement process, it is added to the total response time.

**Followers changing leaders**: Groups may not only exist before, but may also arise during the evacuation [22]. Hereby, leaders with specific properties, such as authority [3] or due to the spatial position [23], may emerge. These emergent groups can be distinguished from social groups with high intragroup social relations, by their steadiness and the attachment among group members and leaders [24]. Phenomenons such as backtracking may be found in social groups [21]. Whereby emergent groups

may only last temporarily, and spatial distances may split group members from the leader [24]. In a dynamic group, a follower may change to a new leader if another leader is closer to the follower [25] and in its visibility [26].

### 3.5    Uncertainty in the Model

Various uncertainties may be encountered in the model and are analyzed in order to receive a robust result regarding the three behaviours. All uncertainties are summarized in Table 1. Encountered values in literature determine the range for the base case. An exact overview may be found at https://github.com/JIRnic.

### 3.6    Verification, Validation and Sensitivity Analysis

Ronchi et al. [27] proposed various verification tests, to verify evacuation models. We applied all verification tests with a positive result.

In order to validate the model, we performed macro and micro validation. Whereby, the micro behavior may be defined as individual behavior of agents [28]. In contrast, macro behavior relates to the overall outcome of the model due to the interaction of agents [28]. For macro validation, we compared the evacuation time to empirical data from Haghani et al. [2]. For Micro validation we contrasted each core behavior to empirical data found in literature or due to face validation [27].

Finally, a higher trust in the built model and increased validity of the model may be achieved with the help of a sensitivity analysis [29]. If the model is sensitive to parameters that also occur in the real world, the trust in the model increases [29]. Sobol sensitivity analysis provides the possibility of conducting a global sensitivity analysis [30]. Further explanation of the Sobol methods can be encountered in [30]. The results indicate that the influence of parameters on the main KPIs is logical and, thus, increases the trust that the suitable model for its purpose was built. All results of tests and the sensitivity analysis can be found on https://github.com/JIRnic.

## 4    Results

Here we represent the model results. We first analyzed the uncertainty space, then we present and discuss the effects of different behaviors.

**Table 1**  Uncertainties in the agent-based model

| Uncertainty | Location of uncertainty | Explanation | Value range in the model |
|---|---|---|---|
| Familiarity | Input data uncertainty | The familiarity may change, depending on the time and location of the building | 0–30 |
| Population | Input data uncertainty | Depending on the building use and time, the population inside the building may change | 100–1200 |
| Percentage groups | Input data uncertainty | Depending on the building and time, the group percentage may differ | 55–70 |
| Group distribution | Input data uncertainty | Different means for a Poisson distribution could be found in the literature | 0.83–1.4 |
| Max crowd density | Input data uncertainty | Different maximum crowd densities can be found in literature | 5–8 |
| Max distance group members | Input data uncertainty | The max distance between group members may vary | 1–6 |
| Min_age | Input data uncertainty | In some areas, no children may be present | 10–20 |
| Max_age | Input data uncertainty | In some areas, no elderly may be available | 65–85 |
| Recognition time distribution | Structural uncertainty | Different recognition time distributions may be found depending on the location and source | Department Store, Restaurant, Office |
| Determination of group leader | Structural uncertainty | In literature, various methods to determine the group leader were encountered | Random, Closest to the exit |

## 4.1   Uncertainty Analysis

In order to investigate how uncertainties in the model influence the overall behavior of the model, we conducted an uncertainty analysis. First we describe the overall behavior in more detail, then feature scoring is utilized.

**Fig. 3** The overall behavior of the base ensemble on the KPIs

**Overall behavior** The overall behavior of the model is summarized in Fig. 3. The plot illuminates the spread of results for each KPI. First, it illustrates that the mean evacuation time for each scenario ranges from 105.72 s until 271.70 s. Whereby the median lies at around 145.58 s. The same emergent behavior may be observed for the 95% percentile of the evacuation time (Fig. 3: left). In contrast, regarding the uncertainties, the mean response time may not be as sensitive as the evacuation times (Fig. 3: middle). However, it shows that the model's max and min response time is highly variable. Finally, we studied intragroup behavior of the model. The mean intragroup distance illustrates a large difference between the scenarios. With an increasing number of group members, the distance between group members grows.

**Feature scoring** The above results demonstrate that uncertainties highly impact the outcome of an evacuation. We utilized feature scoring to identify the relevant influence of uncertainties on the KPIs. A higher score in Fig. 4 indicates a greater influence on the KPI. It shows that the evacuation time is mostly influenced by the population, the familiarity and only minimal from the recognition time distribution. A higher recognition time distribution may lead to longer recognition times and thus impacts the total evacuation time. The same phenomena may be observed for the 95 percentile of the evacuation time. The response time is affected by "the percentage of groups". Whereby the population mainly influences the maximum of the response time.

## 4.2 Leader-Follower Behavior

Fig. 5 compares the three different leader-follower behaviors with the base case without additional leader-follower behaviors. The plot already indicates that backtracking demonstrates a higher evacuation time compared to the base case. In addition, the results from a Mann-Whitney U test confirm this result. Furthermore, the difference between the medians (150.77 for the base case and 202.31 for backtracking) indicates the negative influence on the evacuation time. To conclude, the backtracking of the leader may reduce the speed of groups leading to a higher evacuation time.

The statistical test indicates a difference in the response time between the base case and group gathering. However, the gap between the medians shows only a

| | evacuation100 | evacuation95 | mean response time | SD response time | Max response time | Min response time | Mean intragroup distance total | Mean intragroup distance G2 | Mean intragroup distance G3 | Mean intragroup distance G4 |
|---|---|---|---|---|---|---|---|---|---|---|
| Mean_Poisson_distribution | 0.018 | 0.014 | 0.031 | 0.034 | 0.019 | 0.022 | 0.1 | 0.019 | 0.016 | 0.046 |
| Recognition_distribution_location | 0.14 | 0.12 | 0.0098 | 0.011 | 0.013 | 0.021 | 0.02 | 0.034 | 0.03 | 0.025 |
| determination_of_leader | 0.0095 | 0.009 | 0.0064 | 0.0077 | 0.0088 | 0.012 | 0.014 | 0.015 | 0.012 | 0.015 |
| distance_group_members_setup | 0.016 | 0.015 | 0.013 | 0.014 | 0.016 | 0.022 | 0.43 | 0.42 | 0.59 | 0.39 |
| familiarity_percentage | 0.23 | 0.22 | 0.014 | 0.015 | 0.017 | 0.024 | 0.038 | 0.18 | 0.1 | 0.061 |
| max_age | 0.018 | 0.015 | 0.014 | 0.016 | 0.018 | 0.023 | 0.019 | 0.021 | 0.016 | 0.02 |
| max_congestion | 0.016 | 0.014 | 0.011 | 0.015 | 0.016 | 0.023 | 0.047 | 0.046 | 0.028 | 0.025 |
| min_age | 0.015 | 0.015 | 0.013 | 0.015 | 0.018 | 0.023 | 0.02 | 0.02 | 0.015 | 0.019 |
| percentage_groups | 0.018 | 0.016 | 0.87 | 0.86 | 0.031 | 0.62 | 0.04 | 0.032 | 0.052 | 0.2 |
| population | 0.52 | 0.56 | 0.013 | 0.016 | 0.84 | 0.21 | 0.27 | 0.21 | 0.14 | 0.2 |

**Fig. 4** Feature scoring for the uncertainties. A higher number indicates a greater influence of uncertainties (left side of the figure) on the KPI (bottom of the figure)



**Fig. 5** The evacuation times and response times for each leader-follower behavior compared to the base case

slight divergence. The reason behind the small gap lies in the distribution of the agents. Group members are already located close to each other before the evacuation. This gap increases if the group member is situated further apart. No impact in the evacuation time could be observed by this behavior.

The flexibility of the group increases the evacuation time. Mann-Whitney U test verifies this trend. However, the overall intragroup distance increases. When implementing flexible groups in the model, bigger groups emerge. Overall, these groups demonstrate a higher intragroup distance and lower walking speed, leading to higher evacuation times and increasing the mean distance. All results are summarized in Table 2.

**Table 2**  Results for each leader-follower behavior

| Behavior | KPIs | Mann-Whitney U test: *P*-value | Median base case | Median behavior |
|---|---|---|---|---|
| Backtracking | Evacuation time | < 0.01 | 150.77 | 202.31 |
| | Response time | 0.82 | 27.32 | 27.31 |
| | Mean intragroup distance | < 0.01 | 1.67 | 2.50* |
| | Evacuation time | 0.25 | 150.77 | 152.21 |
| Group gathering | Response time | <0.01 | 27.32 | 27.66* |
| | Mean intragroup distance | 0.10 | 1.67 | 1.63 |
| | Evacuation time | < 0.01 | 150.77 | 195.56* |
| Flexibility of the group | Response time | 0.97 | 27.32 | 27.31 |
| | Mean intragroup distance | < 0.01 | 1.67 | 1.97* |

Significant differences at the *p*-value lower than 0.05 are marked with a *



**Fig. 6**  Results for the multi variant behavior testing. The evacuation times (left) and response time (right) for each combination

## 4.3  Multivariant Behavior Testing

Also, we studied the impact of combinations of leader-follower behaviors. Figure 6 indicates that certain combinations of the behaviors increase the evacuation time compared to the implementation of one leader-follower strategy and the base case. In particular, the combination of group flexibility and backtracking results in the highest increase. An explanation is the higher number of agents per group due to the possibility of changing to another leader. This leads to longer waiting times for the leader as group members may get lost in congestion. Only combinations featuring group gathering have a increased response time. This is logical, as only the group gathering strategy adds to the response time.

# 5 Discussion and Conclusion

Our research question was how do different leader-follower behaviors in groups (backtracking, group gathering, and followers changing leaders) influence the evacuation and response time inside buildings. Therefore, we developed an agent-based model, and utilized an exploratory modeling approach. The main simulation results demonstrated that the group's flexibility and backtracking increased the evacuation time. Whereby group gathering impacts the response time. Additionally, the group distance was monitored and indicated that group gathering reduces the distance between group members during the evacuation. However, backtracking and flexibility of the group displayed diverse results regarding this KPI. It reduced the intragroup distance only for groups with fewer members. However, the reason behind it may be the implementation in the model as evacuees attempt to avoid the patch of other group members and step aside, leading to a higher distance between the group members. For the flexibility of the group, the reason lies in the creation of bigger groups, which generally have a greater intragroup distance [5]. Furthermore, with the help of sampling over the uncertainty space, higher confidence in the results could be achieved. Overall, the results may aid researchers who apply this behavior to understand how different leadership behavior influence the overall evacuation process. Of course, it is essential to remember that a model may never represent the real world, and the outcome is related to the implementation of the behavior in the model. In addition, uncertainties in the model about the input, measurement of the results, agent behavior, and model formalization are present. Furthermore, the model only compared different kinds of leader-follower behavior. Nevertheless, various groups may be observed in an evacuation, with varying decision-making structures [2]. Lastly, no empirical data about different leader-follower behavior are currently available, which increases the difficulties in comparing the model with real-life experiments.

Overall, the results in this study indicated that all additional leader-follower behaviors impact the evacuation performance. Thus, modelers and researchers must include backtracking and group gathering for social groups and flexibility of the group for emergent groups in evacuation models due to their impact on the evacuation performance found in this study. Currently, many models only implement the core leader-follower behavior and neglect the additional behaviors of leader and follower. However, only implementing the core leader-follower behavior in models may lead to wrong conclusions. Policymakers and fire safety engineers may then utilize the models with the included behaviors to prepare buildings for these critical situations and save people's lives.

# References

1. Hogg, M.A., Van Knippenberg, D., Rast, D.E.: European review of social psychology the social identity theory of leadership: theoretical origins, research findings, and conceptual developments. Euro. Rev. Soc. Psychol. **23**(1), 258–304 (2012). https://doi.org/10.1080/10463283.2012.741134
2. Haghani, M., Sarvi, M., Shahhoseini, Z., Boltes, M.: Dynamics of social groups' decision-making in evacuations. Transp. Res. Part C Emerg. Technol. **104**, 135–157 (2019). https://doi.org/10.1016/j.trc.2019.04.029
3. Jones, B.K., Hewitt, J.A.: Leadership and Group Formation in High-Rise Building Evacuations, pp. 513–522 (1986). https://doi.org/10.3801/iafss.fss.1-513
4. van der Wal, C.N., Robinson, M.A., Bruine de Bruin, W., Gwynne, S.: Evacuation behaviors and emergency communications: An analysis of real-world incident videos. Safety Sci. **136**, 105121 (2021). https://doi.org/10.1016/J.SSCI.2020.105121
5. Moussaïd, M., Perozo, N., Garnier, S., Helbing, D., Theraulaz, G.: The walking behaviour of pedestrian social groups and its impact on crowd dynamics. PLoS One **5**(4), 10047 (2010). https://doi.org/10.1371/journal.pone.0010047
6. Bernardini, G., Ciabattoni, L., Quagliarini, E., D'Orazio, M.: Cognitive buildings for increasing elderly fire safety in public buildings: design and first evaluation of a low-impact dynamic wayfinding system. Lecture Notes Electr. Eng. **725**, 101–119 (2021). https://doi.org/10.1007/978-3-030-63107-9_8
7. Cuesta, A., Abreu, O., Alvear, D.: Methods for measuring collective behaviour in evacuees. Safety Sci. **88**, 54–63 (2016). https://doi.org/10.1016/J.SSCI.2016.04.021
8. Li, J., Xue, B., Wang, D., Xiao, Q.: Study on a new simulation model of evacuation behavior of heterogeneous social small group in public buildings. J. Appl. Sci. Eng. **24**(4), 467–475 (2021). https://doi.org/10.6180/jase.202108_24(4).0002
9. Wang, J., Nan, L., Lei, Z.: Small group behaviors and their impacts on pedestrian evacuation. In: Proceedings of the 2015 27th Chinese Control and Decision Conference, CCDC 2015, pp. 232–237. Institute of Electrical and Electronics Engineers Inc. (2015). https://doi.org/10.1109/CCDC.2015.7161696
10. Zheng, X., Zhong, T., Liu, M.: Modeling crowd evacuation of a building based on seven methodological approaches. Build. Environ. **44**(3), 437–445 (2009). https://doi.org/10.1016/j.buildenv.2008.04.002
11. Jennings, N.R.: On agent-based software engineering. Artif. Intell. **117**, 277–296 (2000)
12. Ronchi, E., Reneke, P.A., Peacock, R.D.: A method for the analysis of behavioural uncertainty in evacuation modelling. Fire Technol. **50**(6), 1545–1571 (2014). https://doi.org/10.1007/s10694-013-0352-7
13. Weaver, C.P., Lempert, R.J., Brown, C., Hall, J.A., Revell, D., Sarewitz, D.: Improving the contribution of climate model information to decision making: the value and demands of robust decision frameworks. Wiley Interdiscip. Rev. Clim. Change **4**(1), 39–60 (2013). https://doi.org/10.1002/WCC.202
14. Kwakkel, J.H.: The exploratory modeling workbench: an open source toolkit for exploratory modeling, scenario discovery, and (multi-objective) robust decision making. Environ. Modell. Softw. **96**, 239–250 (2017). https://doi.org/10.1016/J.ENVSOFT.2017.06.054
15. Auping, W.L.: Modelling Uncertainty: Developing and Using Simulation Models for Exploring the Consequences of Deep Uncertainty in Complex Problems. Ph.D. thesis, Delft University of Technology (2018). https://doi.org/10.4233/UUID:0E0DA51A-E2C9-4AA0-80CC-D930B685FC53
16. Yang, D.Y., Frangopol, D.M.: Risk-based vulnerability analysis of deteriorating coastal bridges under hurricanes considering deep uncertainty of climatic and socioeconomic changes. ASCE-ASME J. Risk Uncertainty Eng. Syst. Part A Civ. Eng. **6**(3), 04020032 (2020). https://doi.org/10.1061/AJRUA6.0001075

17. ...Grimm, V., Berger, U., Bastiansen, F., Eliassen, S., Ginot, V., Giske, J., Goss-Custard, J., Grand, T., Heinz, S.K., Huse, G., Huth, A., Jepsen, J.U., Jørgensen, C., Mooij, W.M., Müller, B., Pe'er, G., Piou, C., Railsback, S.F., Robbins, A.M., Robbins, M.M., Rossmanith, E., Rüger, N., Strand, E., Souissi, S., Stillman, R.A., Vabø, R., Visser, U., DeAngelis, D.L.: A standard protocol for describing individual-based and agent-based models. Ecol. Modell. **198**(1–2), 115–126 (2006). https://doi.org/10.1016/J.ECOLMODEL.2006.04.023

18. Grimm, V., Berger, U., DeAngelis, D.L., Polhill, J.G., Giske, J., Railsback, S.F.: The ODD protocol: a review and first update. Ecol. Modell. **221**(23), 2760–2768 (2010). https://doi.org/10.1016/J.ECOLMODEL.2010.08.019

19. Ronchi, E.: Developing and validating evacuation models for fire safety engineering. Fire Safety J. **120**, 103020 (2021). https://doi.org/10.1016/J.FIRESAF.2020.103020

20. Köster, G., Treml, F., Seitz, M., Klein, W.: Validation of crowd models including social groups. Pedestrian Evacuation Dyn. **2012**, 1051–1063 (2014). https://doi.org/10.1007/978-3-319-02447-9_87

21. Lu, L., Chan, C.Y., Wang, J., Wang, W.: A study of pedestrian group behaviors in crowd evacuation based on an extended floor field cellular automaton model. Transp. Res. Part C Emerg. Technol. **81**, 317–329 (2017). https://doi.org/10.1016/j.trc.2016.08.018

22. Quarantelli, E.L.: Emergent behaviors and groups in the crisis time of disasters. In: Individuality and Social Control: Essays in Honor of Tamotsu Shibutani, pp. 47–68 (1995)

23. Lombardi, M., Warren, W.H., di Bernardo, M.: Nonverbal leadership emergence in walking groups. Sci. Rep. **10**(1) (2020). https://doi.org/10.1038/S41598-020-75551-2

24. Fang, J., El-Tawil, S., Aguirre, B.: Leader-follower model for agent based simulation of social collective behavior during egress. Safety Sci. **83**, 40–47 (2016). https://doi.org/10.1016/J.SSCI.2015.11.015

25. Ji, Q., Gao, C.: Simulating crowd evacuation with a leader-follower model. IJCSES Int. J. Comput. Sci. Eng. Syst. **1**(4) (2007)

26. Mao, Y., Fan, X., Fan, Z., He, W.: Modeling group structures with emotion in crowd evacuation. IEEE Access **7**, 140010–140021 (2019). https://doi.org/10.1109/ACCESS.2019.2943603

27. Ronchi, E., Kuligowski, E.D., Nilsson, D., Peacock, R.D., Reneke, P.A.: Assessing the verification and validation of building fire evacuation models. Fire Technol. **52**(1), 197–219 (2016). https://doi.org/10.1007/s10694-014-0432-3

28. Moss, S., Edmonds, B.: Sociology and simulation: statistical and qualitative cross-validation. Am. J. Sociol. **110**(4), 1095–1131 (2005). https://doi.org/10.1086/427320

29. Smith, E.D., Szidarovszky, F., Karnavas, W.J., Bahill, A.T.: Sensitivity analysis, a powerful system validation technique. Open Cybern. Systemics J. **2**, 39–56 (2008)

30. Sobol, I.M.: Global sensitivity indices for nonlinear mathematical models and their Monte Carlo estimates. Math. Comput. Simul. **55**(1–3), 271–280 (2001). https://doi.org/10.1016/S0378-4754(00)00270-6

# Relation Between the Public and the Private and Evolution of Food Sharing



**Elpida Tzafestas**

**Abstract**  We are studying food sharing and the conditions of its evolution in sedentary, agricultural societies. We start from the conjecture that when food, and by extension any type of wealth, can be stored for public use in greater quantities than what can be stored for private use, then sharing evolves, otherwise it does not. We present a simulated environment where agents represent families that produce, consume and store food in private or in public stores. The relative capacities of the private and public stores are responsible for the evolution or not of generalized sharing in the population. Although storage capacity is represented in the model as a technological parameter, it is not purely technological, but it is also the result of social and cultural choices and it goes hand in hand with communal or individualist behavioral profiles and stances. We show that when the public sphere is given a culturally greater value than the private sphere and therefore public storage capacity is higher than private capacity, then sharing emerges. This is more pronounced in risky and unstable environments. Inversely, if private capacity rises to very high levels compared to the public one, sharing collapses. Similar results are obtained even when other processes are present, such as costly sharing and theft. Our conclusion is therefore that sharing can emerge in a sedentary population when storage is possible and it will be maintained as long as individuals can store only limited amounts of food or other wealth. We delineate the implications of our study for further research on sharing and inequality.

**Keywords**  Sharing · Food sharing · Evolution · Storage · Theft

## 1   Introduction

Food transfer and sharing is very common among humans and especially pervasive in small-scale populations and societies, namely foragers or agricultural-foragers. Sharing is therefore one important field of study in anthropology, behavioral ecology,

E. Tzafestas (✉)
Laboratory of Cognitive Science, Department of History and Philosophy of Science, National and Kapodistrian University of Athens, University Campus, 15771 Ano Ilisia, Athens, Greece
e-mail: etzafestas@phs.uoa.gr

human evolution and other disciplines [1, 2]. One fundamental feature of sharing is that it represents unresisted transfer between unrelated individuals and that it is inextricably connected to social living and bonding. The evolutionary origins of sharing and its trajectory from other animal and especially primate behaviors are a vivid object of research and extensive theorizing. The most prominent theories of evolution of human sharing include: reciprocal altruism (either through aggressive sharing and assertive reciprocation [3] or through explicit reciprocity [4]), group cooperation (where groups of sharers have a selective advantage over groups of non-sharers [5–7]), direct reciprocity or risk reduction (where giving now allows to ask later if resources are highly volatile [8, 9]), tolerated theft (where tolerance is preferred to fight [10]), costly signaling (where sharing serves to show-off and to demonstrate reproductive quality [11]) and processual approaches (where need and behavior are shared and thus food as well [12, 13]). The relation of sharing with other social processes and phenomena is also a permanent object of study [14, 15].

Although sharing is a cooperative behavior, it is not altruistic [3], because (a) it does not always entail a cost to the donor while altruistic helping always does, and (b) it may be enforced through social norm or direct coercion whereas it would appear difficult to enforce altruism. Thus sharing is a controlled and reasoned process, despite the fact that its specific forms are internalized and automated by immersion in a specific cultural context.

We adopt at first a direct reciprocity or risk-reduction approach [8, 9, 16, 17] with an added feature of food storage, since if excess food cannot be stored then it has to either be shared or thrown away. Thus we expect storage in highly unstable environments and storage constraints in general to drive evolution of sharing and to promote food sharing [18, 19]. Although risk reduction has been studied essentially for hunter-gatherer and forager populations, where large chunks of food, such as game meat, become available only very sparsely [20–22], we are interested in sedentary agricultural populations. This has the disadvantage that some sharing practice and culture may already be present as a heritage of earlier traditions from more primitive, forager, populations, but it allows us to explicitly model and study the relation of sharing with surplus and storage and to later theorize about how this can lead to higher complexity of sociopolitical structures and wealth inheritance [23, 24]. We maintain an eye on tolerated theft, however, because this is very common in higher mammal and especially primate species and so we suspect that it may serve as a base on which explicit sharing could evolve in humans. Some modeling efforts have investigated various aspects of sharing and its relation to various processes [25–28]. Our own focus is on the conditions for evolution of sharing.

This paper is structured as follows. In Sect. 2, we present the agent and environment model and in Sect. 3 we present simulation results for the base condition (unlimited storage). Next, in Sect. 4 we report the results of experiments with the basic model that show that the relative values of public and private storage capacity define the evolutionary dynamics of the system. In Sects. 5–7 we investigate what happens when sharing has a cost, or when the possibility of theft is introduced or

when agents are allowed to decide independently to give to or to receive from the public store. In Sect. 8 we study "shocks", where the storage capacities are abruptly reinitialized, to find out that evolutionary dynamics adapt accordingly. Final thoughts are given in the last section.

## 2 The Model

We are studying a fairly simple agent-based situated production model that pertains directly to human preindustrial agricultural populations, although it may be argued to apply to other forms of primitive economy as well, provided there is no trade. The model uses a spatial grid of $30 \times 30$ cells where a number of agents representing human families/groups live, produce and occasionally move. In what follows, we use the terms agent and family interchangeably. Every cell has a level of fertility and each family has a level of technological ability that allows it to extract food. The environment has a degree of instability, i.e., a probability with which a percentage of its production is lost (for climatic or natural resource reasons). Food that is not immediately consumed is stored for the future in private family stores or in public ones and the families decide about whether to store publicly and how much. All stores may be of limited or unlimited capacity.

Families grow or shrink with a constant birth and death rate, respectively. When the size of a family exceeds a size threshold, the family splits to two and the newborn family takes half the people and half the stored food. If the fertility level of the current cell is insufficient, the newborn family migrates. If a family cannot sustain itself within the current cell either because fertility is insufficient or because its stored food has fallen below a security level, it migrates as well. If a migration cell with sufficient fertility cannot be found, a family may initiate a war against a neighbour. Such a survival war is initiated against the richest neighbour and the aggressor is supposed to be always successful. The aggressor and winner then steals all the privately stored food of the victim as well as its share in the public store. The victim migrates elsewhere or dies if this is not possible. The general algorithm is given below:

1. Production locally according to size of family, technological means and place productivity
2. Public sharing (if applicable): a proportion of the production (= sharing rate * production) is sent to the public store in the agent's position
3. surplus = rest of production (after sharing)—current need
4. If surplus > 0, then store it privately
5. Else consume (-surplus) from publicly stored food in its position or in neighbouring positions

6. If still in need, launch survival war against the richest neighbour
     If no such neighbour exists, try to migrate in a rich nearby place
     If nothing works, die (starve)

Table 1 presents the most important parameters of the model. Parameter diversity allows the emergence of population differences.

All the agent parameters are inherited during family split and in most of the basic experiments (Sects. 3 and 4) there is a 5% mutation rate (probability to flip the sharing gene and reset the sharing rate). This basic model leads to Malthusian population evolution, where after a while agents fill the whole 2D-array of cells and the population size stabilizes around a value that may be regarded as the carrying capacity of the environment for certain instability characteristics and technological abilities (see Fig. 1). The environmental and behavioral parameters have been tuned to allow the population to stabilize fairly quickly (in 1000 to 2000 cycles) to this Malthusian state.

The parameters have also been tuned to demonstrate the evolutionary trends in terms of sharing. In the experiments that follow we initialize each family in the population with the sharing gene ON with 50% probability and with a sharing rate uniformly distributed between 0 and 0.5 (thus on average 0.25). This setup allows us to quickly explore the evolutionary trends of the system without waiting too long for the pro- or anti- sharing mutations to spread from scratch. The evolutionary dynamics of sharing are generally of one of two forms: either full sharing or no sharing emerges

**Table 1** Parameters of the basic model

| *Environment* | | | |
|---|---|---|---|
| Size | $30 \times 30$ | No. of agents at t = 0 | N = 100 |
| Environmental instability | 0 or 0.2 or 0.5 | Cell fertility | Uniform (400–600) |
| Environmental loss of production rate | 0 or 0.5 or 0.8 | Max. storage capacity/cell | 10 to 1000 or unlimited ($-$ 1) |
| *Agent behavior* | | | |
| Migration cost | 1 or 2 | Technology | Uniform (1–5) |
| Need for food | 10 or 11 | Maximum family size | 40 or 41 |
| Vision | 1 to 4 cells | Max. storage capacity | 10 to 1000 or unlimited ($-$ 1) |
| War vision | 1 to 4 cells | Sharing gene | On/Off (with initial probability 50%, 50%) |
| Food security level | 1 or 2 | Sharing rate | Uniform (0–0.5) |
| | | Mutation rate | 0.05 |
| Birth rate | 0.3 | Death rate | 0.2 |

**Fig. 1** (x: time, y: number of live families in 1000 s) Typical outcome for initial N = 200, technology = 1 for all agents. The population stabilizes around 1400 families that fill the 30 × 30 cell grid (1 or 2 families per cell). The population stabilizes to the Malthusian limit independently of the initial number of families. The speed of convergence to the limit may differ according to the various behavioral parameters



**Fig. 2** (x: time, y: avg. sharing gene value in the population) Two typical sharing gene evolutionary movements. Left: A population where sharing disappears. Right: A population where sharing is established

(see Fig. 2). Intermediate degrees of sharing are very rare and they are found almost exclusively in cases where the population has not stabilized within the experimental timeframe used (3000 cycles unless otherwise stated). The presence of mutations has not been found to influence these trends, so most of the advanced behavioral experiments (Sects. 5–8) have been carried out without mutation.

## 3 Reference condition

We have run experiments to evaluate whether sharing evolves and to what degree in various environmental conditions and technological setups. In all cases, we start from 50% presence of the sharing gene in the initial population with an average rate of 25% (0.25, uniformly distributed between 0 and 0.5) and we measure and compare

the average sharing gene and sharing rate in the final society after stabilization to a Malthusian condition. All results of experiments in this and the following sections are averages of 20 runs. This value has been found empirically to ensure that the resulting standard deviations of all experiment metrics are very low compared to averages.

| | No instability | | High instability | |
|---|---|---|---|---|
| | Tech 1 uniform | Tech 5 diverse | Tech 1 uniform | Tech 5 diverse |
| Live | 1481.86 | 4722.72 | 718.96 | 3016.62 |
| Sharing gene | 0.365 | 0.131 | 0.82 | 0.347 |
| Sharing rate | 0.023 | 0.029 | 0.053 | 0.069 |

In the above reference experiment, we give the results of basic productive versus more advanced agents (technology = 1 uniformly in the population or technology uniformly distributed between 1 and 5) in stable (instability = 0) or extremely unstable environments (instability = 0.5, loss of production rate = 0.8). Here, public and private storage capacities are unlimited.

We observe that sharing is far more common and with higher rates in unstable environments and for less productive agents. The numbers of live families depend on both these parameters (technological productivity and environmental instability) and are around the Malthusian limit in all cases. According to the previous section, the sharing gene results ought to be read as proportions of runs that have led to full sharing. We obtain similar results if we disable survival wars or if there is diversity in terms of a "survival war gene" and for other variations of the original model.

## 4   Storage and Sharing

We now proceed to test our basic hypothesis that low private combined with high public storage capacities favor the evolution of sharing, whereas high private combined with low public storage capacities favor selfishness. Storage capacity is both a technological and a social and cultural parameter. The technological dimension is obvious: it has to do with how much food one can store (one can store more meat if smoking is possible, and one can store more fruit if a refrigerator is available, etc.). The social and cultural dimensions are less obvious. If community life is of high value, more effort and investment will be put into creating and maintaining a high capacity public store and less into high capacity private stores. High capacity stores can also be hard to defend against raiders and thieves, so that a high capacity private store will be much harder to defend individually than an equally high capacity public store that can be defended collectively (with shifts of guards and the like). Violation costs may also be much higher for public offenders than for private ones; as a consequence, storage capacity can be higher for public than for private stores.

All these parameters taken together are represented as a lumped storage capacity parameter.

We run populations in the environmental and technological conditions of the previous section (technology = 1 uniformly or 1–5 randomly and stable or highly unstable environment). We use four different private–public storage combinations as given in the following table. With a base public capacity of 20, we test the cases of comparatively low private capacity (50), medium capacity (500) and high capacity (1000). We also test the case of diverse private storage capacities (from 50 to 1000 for every agent), thus low for some agents and high for others. As before, we measure average sharing gene and sharing rate as well as average private storage capacity in the final population. The results show that, in all cases, a fairly low private storage capacity leads to almost full sharing evolution (close to 100%), while an extremely high private storage capacity leads to almost full sharing collapse (0 or close to 0). Intermediate values of private storage capacity (such as 500) lead to intermediate degrees of sharing presence in the population, while in diverse populations with private storage capacities randomly drawn from 50 to 1000 the lower private storage capacity agents are quickly wiped out by evolution, the average private storage capacity rises and sharing sinks or even disappears.

Note also that as average storage capacity rises, so does the population size, because there is more stored food available to allow survival in case of temporal fluctuations. But, in all cases and for the same reason, the populations at the Malthusian limit are of a lower size than those for unlimited storage capacities of the previous section.

| $PU = 20$ | No instability | | High instability | |
|---|---|---|---|---|
| | Tech 1 uniform | Tech 5 diverse | Tech 1 uniform | Tech 5 diverse |
| *PRI = 50* | | | | |
| Live | 996.5 | 3613.9 | 458.58 | 2192.82 |
| Sharing gene | 1.0 | 0.98 | 1.0 | 1.0 |
| Sharing rate | 0.251 | 0.243 | 0.264 | 0.246 |
| *PRI = 500* | | | | |
| Live | 1160.18 | 3696.26 | 545.64 | 2263.92 |
| Sharing gene | 0.641 | 0.2 | 0.267 | 0.157 |
| Sharing rate | 0.205 | 0.143 | 0.092 | 0.08 |
| *PRI = 1000* | | | | |
| Live | 1234.74 | 3879.24 | 698.48 | 2800.14 |
| Sharing gene | 0 | 0.005 | 0 | 0.023 |
| Sharing rate | 0 | 0.054 | 0 | 0.089 |
| *PRI = 50–1000* | | | | |
| Live | 1178.52 | 3543.6 | 629.3 | 2530.88 |
| Sharing gene | 0.001 | 0.046 | 0 | 0.082 |
| Sharing rate | 0.008 | 0.047 | 0 | 0.098 |

(continued)

| PU = 20 | No instability | | High instability | |
|---|---|---|---|---|
| | Tech 1 uniform | Tech 5 diverse | Tech 1 uniform | Tech 5 diverse |
| Storage | 889.258 | 869.783 | 814.462 | 790.368 |

## 5   Costs and Sharing

The private and public storage capacities are therefore responsible for the evolution or collapse of sharing in a population, but sharing is not always free lunch and often comes with a cost. More specifically, sharing is usually associated with elevated degrees of bonding within a population, with extensive community life participation, religious rituals and so on [1, 5, 29]. These activities have a cost, at least because they demand time that is to the detriment of other more directly productive activities. The size of the cost can be crucial: a low sharing cost may boost sharing if the expected profit from sharing is high, and, vice versa, a high sharing cost may dissolve sharing bonds in a stressful environment of production where resources are very limited for other reasons.

| | Tech 1 uniform, no instability | | | Tech 5 diverse, high instability | | |
|---|---|---|---|---|---|---|
| PU = 300, PRI = 50 | Cost = 10 | Cost = 45 | Cost = 10–45 | Cost = 10 | Cost = 45 | Cost = 10–45 |
| Live | 1105.78 | 950.44 | 1160.22 | 2610.6 | 2058.46 | 2473.8 |
| Sharing gene | 1.0 | 0 | 1.0 | 0.986 | 0 | 0.983 |
| Sharing rate | 0.083 | 0 | 0.137 | 0.237 | 0 | 0.255 |
| Sharing cost | 10 | 45 | 15.716 | 10 | 45 | 19.524 |

We report above indicative results of an experiment where an otherwise sharing-friendly (or selfishness-hostile) environment of high public and low private storage capacity may see a collapse of sharing because of high cost (such as 45 per time step), while for low costs (such as 10 per time step) sharing stills evolves. As before, we measure average sharing gene and sharing rate as well as average sharing cost in the final population. It is noteworthy that in populations with diversified degrees of sharing cost, that correspond to families either "willing" or "reluctant" to share (and hence low-cost or high-cost, respectively), the low-cost, sharing-prone families overtake the population (as seen in the average sharing cost of the final population, that is below the initial average of 28) and sharing again evolves.

## 6 Theft and Sharing

In human, as well as in animal populations, cooperation is not always full and cheaters, kleptoparasites and similar behavioral profiles are very common [30, 31]. We test whether a degree of kleptoparasitism in the form of theft can be maintained in a sharing population or whether sharing collapses right away in this case. We define an additional **"public theft" gene** that makes an agent a potential thief: such an agent, decides with a certain probability (**theft rate**) to steal from the public store instead of producing food. The rest of its behavior remains the same: it may share food as usual, and may ask and obtain food from the public store if still in need. Thus, a sharing thief agent will occasionally steal instead of producing and it will participate in the rest of the activities as all others.

| | Tech 1 uniform (public theft) | | Tech 5 diverse (public theft) | | Tech 1 uniform (public theft) | | | |
|---|---|---|---|---|---|---|---|---|
| | *Unlimited* storage | | *Unlimited* storage | | *PU = 20, PRI = 20* | | *PU = 20, PRI = 1000* | |
| | No | High | No | High | No | High | No | High |
| Live | 1474.84 | 716.54 | 4254.0 | 2879.54 | 986.76 | 452.58 | 1239.4 | 701.82 |
| Sh. gene | 0.482 | 0.786 | 0.118 | 0.419 | 1.0 | 1.0 | 0 | 0 |
| Sh. rate | 0.033 | 0.05 | 0.019 | 0.105 | 0.238 | 0.247 | 0 | 0 |
| Thieves (rate) | 0.588 (0.468) | 0.514 (0.335) | 0.066 (0.039) | 0.017 (0) | 0.577 (0.469) | 0.488 (0.478) | 0.539 (0.448) | 0.417 (0.433) |
| Sharer-thieves | 0.285 | 0.412 | 0 | 0.017 | 0.577 | 0.488 | 0 | 0 |

In the above experiment, we initialize all agents with 50% theft gene and average theft rate uniformly distributed between 0 and 1 (thus 0.5 on average) and we measure average sharing gene and sharing rate as well as average theft gene, average theft rate and average proportion of agents that are both thieves and sharers in the final population. The results show that (a) in the case of low technological ability with unlimited public and storage capacities a high proportion of thieves with equally high theft rates remain in the population, but this proportion plummets for more advanced technological abilities, and (b) in the case of limited storage capacities the direction of evolution of sharing (toward full sharing or extinction of sharing) is not affected by the presence of thieves and a very high thief proportion with equally high theft rates remains in the population. For more advanced technological abilities and limited storage capacities, the proportion of thieves plummets as before, again without affecting the ultimate direction of the evolution of sharing (results omitted for lack of space). Other types of theft have been also examined, most notably pairwise repeated theft, where the thief repeatedly chooses neighbours to steal from as much as possible until its needs are satisfied, and the results have been more or less similar.

Thus we find that the evolution of sharing is not affected by kleptoparasitism, which on the contrary may be wiped out by sharing. But the issue of theft and klepto-parasitism deserves a thorough study in itself, especially because tolerated theft is

one of the theories behind evolution of sharing and predates direct and explicit human sharing, as has been found in studies with animals and especially primates [10].

## 7  Split Sharing

We also examine the case where an agent may be independently a sharer (giver) and/or a receiver of shared food from the public store. In the following experiment, we initialize all agents in the population as givers with 50% probability and as receivers with 50% probability (thus we have initially approximately 25% givers, 25% receivers, 25% full sharers, i.e., simultaneously givers and receivers, and 25% "asocial" agents, i.e., neither givers nor receivers). As before, we measure average sharing gene and sharing rate as well as average proportions of givers, receivers, full sharers and average private storage capacity in the final population.

| PU = 20 | Tech 1 uniform, no instability, public split sharing | | | Tech 5 diverse, high instability, public split sharing | | |
|---|---|---|---|---|---|---|
| PRI | 20 | 1000 | 20–1000 | 20 | 1000 | 20–1000 |
| Live | 973.54 | 1315.92 | 1250.44 | 2150.7 | 2904.32 | 2487.72 |
| Givers | 0.842 | 0.387 | 0.47 | 0.511 | 0.381 | 0.559 |
| Sharing rate | 0.125 | 0.134 | 0.14 | 0.102 | 0.141 | 0.179 |
| Receivers | 1.0 | 0.003 | 0.039 | 0.922 | 0.04 | 0.132 |
| Full sharers | 0.842 | 0 | 0.034 | 0.506 | 0 | 0.038 |
| Final storage | 20 | 1000 | 899.019 | 20 | 1000 | 827.424 |

The results for the "public split sharing" scheme presented above parallel our former results in that the proportion of givers rises to a very high value and the proportion of receivers comes close to 100% for the cases that have been found before to lead to evolution of sharing (low private storage capacity). In turn, the proportion of receivers gets close to 0 for high private storage capacity, although interestingly the proportion of givers remains at intermediate values (over 30%) which might act as a safety valve for the future. As is shown in the table, most of the receivers are also givers (they are thus full sharers), whereas most of the givers are not receivers, because the givers are very often many more than the full sharers.

Thus we have found that sharing evolves in the same conditions as before, even if its actual form is a different "split" form that has independent components for giving and receiving. We have also experimented with alternate forms of giving/receiving, and more specifically with "private" sharing where an agent may share and give/receive privately with others. Similar results are also obtained in these cases. Both these issues of split sharing and private sharing are important and would deserve specialized studies, because both these mechanisms appear to ask for lower cognitive load and less cooperative effort than full public sharing. As such, these mechanisms may hint into the deeper prosocial mechanisms of sharing and exchange that are

found in pre-agricultural and hunter-gatherer societies that are probably the first step toward full-fledged deliberate and complex sharing forms.

## 8  Shocks

We finally examine what happens in case of a "shock", i.e. of an abrupt reinitialization of the storage capacities, from low to high or from high to low. The first corresponds to a technological revolution while the second might be the result of an environmental disruption or of a new social, cultural or political order that for example puts an upper limit to private property size. In the following experiments, we run the system as before for 3000 cycles in the first condition (Part I), then we reinitialize just the storage parameters (without affecting the populations and the rest of their behavioral parameters) and we rerun for another 3000 cycles (Part II). The results show that the proportion of sharing agents restabilizes to the predicted value of the new condition (the same goes for the size of the population, that regains a new Malthusian limit). The rate of sharing does not change dramatically and remains around its expected value of 0.25. Thus the tendency of the population concerning sharing is extremely strong and is maintained even if the system starts from an altogether different initial condition.

|              | Tech 1 uniform, no instability |                        | Tech 1 uniform, high instability |                        |
|--------------|--------------------------------|------------------------|----------------------------------|------------------------|
| Storage      | Part I: Low (50)               | Part II: High (1000)   | Part I: Low (50)                 | Part II: High (1000)   |
| Live         | 989.55                         | 1298.2                 | 454.25                           | 697.9                  |
| Sharing gene | 0.995                          | 0.017                  | 0.995                            | 0.015                  |
| Sharing rate | 0.244                          | 0.258                  | 0.266                            | 0.267                  |

## 9  Discussion

We have shown that evolution of sharing is controlled by the relation between private and public storage capacities, where capacities are complex variables representing technological, social and cultural factors related to storage. More specifically, high public and low private storage capacities lead to evolution of full sharing in the population, whereas low public and high private storage capacities lead to sharing extinction. This feature is maintained if additional intricacies are present in the sharing environment (sharing costs, kleptoparasitism, independent givers and receivers) and if the population suffers storage reinitialization shocks. Because sharing is favored by more unstable environments under rudimentary technological conditions, we believe that we found a possible evolutionary pathway for emergence and stabilization of

sharing in harsh conditions and for disappearance of sharing later when storage technology becomes better—which in turn allows individualism to grow and more effort to be put in private storage buildup and maintenance.

There are at least two issues that deserve further investigation. Firstly, do our results about theft not affecting final sharing levels mean that preexisting theft and tolerated theft may give rise evolutionarily to explicit sharing, as may be deduced from some anthropological studies [2]? This is to be tested in a sequel model. Secondly, how does sharing relate with inequality? Does extensive sharing mean low inequality or does high inequality lead to haves and have-nots becoming frequent givers and receivers respectively, thus creating an unequal "class-based" society with charity toward the poor, as an initial reading of our results of Sect. 7 gives? We can think of other secondary issues as well that merit further examination, such as whether kinship-based sharing can evolve or the effect of leaders and social role models. In any case, our results open up many independent routes for further study.

# References

1. Gurven, M.: To give and to give not: the behavioral ecology of human food transfers. Behav. Brain Sci. **27**(4), 543–583 (2004)
2. Kaplan, H., Gurven, M.: The natural history of human food sharing and cooperation: a review and a new multi-individual approach to the negotiation of norms. In: Gintis, H., et al. (eds.) Moral Sentiments and Material Interests, The Foundations of Cooperation in Economic Life, Ch. 3, pp. 75–113 (2005)
3. Moore, J.: The evolution of reciprocal sharing. Ethol. Sociobiol. **5**(1), 5–14 (1984)
4. Jaeggi, A.V., Gurven, M.: Natural cooperators: food sharing in humans and other primates. Evol. Anthropol. **22**(4), 186–195 (2013)
5. Nettle, D., Panchanathan, K., Shakti Rai, T., Fiske, A.P.: The evolution of giving, sharing and lotteries. Curr. Anthropol. **52**(5), 747–756 (2011)
6. Boyd, R.: The evolution of reciprocity in sizable groups. J. Theor. Biol. **132**(3), 337–356 (1988)
7. Winterhalder, B.: Social foraging and the behavioral ecology of intragroup resource transfers. Evol. Anthropol. **5**(2), 46–57 (1996)
8. Kaplan, H.S., Schniter, E., Smith, V.L., Wilson, B.J.: Risk and the evolution of human exchange. Proc. R. Soc. Lond. B Biol. Sci. **279**(1740), 2930–2935 (2012)
9. Kaplan, H.S., Schniter, E., Smith, V.L., Wilson, B.J.: Experimental tests of the tolerated theft and risk-reduction theories of resource exchange. Nat. Human Behav. **2**(6), 383–388 (2018)
10. Blurton-Jones, N.G.: Tolerated theft, suggestions about the ecology and evolution of sharing, hoarding and scrounging. Social Sci. Inform. **26**(1), 31–54 (1987)
11. Bliege Bird, R., Smith, E.A., Bird, D.W.: The hunting handicap: costly signaling in male foraging strategies. Behav. Ecol. Sociobiol. **50**(1), 9–19 (2001)
12. Smith, D., Dyble, M., Major, K., Page, A.E., Chaudhary, N., Salali, G.D., Thompson, J., Vinicius, L., Bamberg Migliano, A., Mace, R.: A friend in need is a friend indeed: need-based sharing, rather than cooperative assortment, predicts experimental resource transfers among Agta hunter-gatherers. Evol. Hum. Behav. **40**(1), 82–89 (2019)
13. Widlok, T.: Extending and limiting selves: a processual theory of sharing. In: Lavi, N., Friesem, D.E. (eds.) Towards a Broader View of Hunter-Gatherer Sharing, Ch. 2, pp. 25–38 (2019)
14. Lee, R.: Sociality, selection, and survival: simulated evolution of mortality with intergenerational transfers and food sharing. Proc. Natl. Acad. Sci. **105**(20), 7124–7128 (2008)
15. Pickles, A.J.: Transfers—a deductive approach to gifts, gambles, and economy at large. Curr. Anthropol. **61**(1), 11–29 (2020)

16. Hao, Y., Armbruster, D., Cronk, L., Aktipis, C.A.: Need-based transfers on a network: a model of risk-pooling in ecologically volatile environments. Evol. Hum. Behav. **36**(4), 265–273 (2015)
17. Aktipis, C.A., Cronk, L., de Aguiar, R.: Risk-pooling and herd survival: an agent-based model of a maasai gift-giving system. Hum. Ecol. **39**(2), 131–140 (2011)
18. Ringen, E.J., Duda, P., Jaeggi, A.V.: The evolution of daily food sharing: a Bayesian phylogenetic analysis. Evol. Hum. Behav. **40**(4), 375–384 (2019)
19. Ember, C.R., Skoggard, I., Ringen, E.J., Farrer, M.: Our better nature: does resource stress predict beyond-household sharing. Evol. Hum. Behav. **39**(4), 380–391 (2018)
20. Dowling, J.H.: Individual ownership and the sharing of game in hunting societies. Am. Anthropol. **70**(3), 502–507 (1968)
21. Hawkes, K., O'Connell, J.F., Blurton Jones, N.G.: Hunting and nuclear families—some lessons from the Hadza about men's work. Curr. Anthropol. **42**(5), 681–709 (2001)
22. Knight, B.: The anonymity of the hunt—a critique of hunting as sharing. Curr. Anthropol. **53**(3), 334–355 (2012)
23. Mattison, S.M., Smith, E.A., Shenk, M.K., Cochrane, E.E.: The evolution of inequality. Evol. Anthropol. **25**(4), 184–199 (2016)
24. Hayden, B.: Archaeological pitfalls of storage. Curr. Anthropol. **61**(6), 763–793 (2020)
25. Thébaud, O., Locatelli, B.: Modelling the emergence of resource-sharing conventions: an agent-based approach. J. Artif. Soc. Social Simul. 4(2) (2001). https://www.jasss.org/4/2/3.html
26. Younger, S.M.: Discrete agent simulations of the effect of simple social structures on the benefits of resource sharing. J. Artif. Soc. Social Simul. 6(3) (2003). https://www.jasss.org/6/3/1.html
27. Bots, P.W.G., Barreteau, O., Abrami, G.: Measuring solidarity in agent-based models of resource sharing situations. Adv. Complex Syst. **11**(2), 337–356 (2008)
28. Baker, M.J., Swope, K.J.: Sharing, gift-giving, and optimal resource use incentives in hunter-gatherer society. Econ. Govern. **22**(2), 119–138 (2021)
29. Gintis, H., Bowles, S., Boyd, R., Fehr, E.: Moral sentiments and material interests: origins, evidence and consequences. In: Gintis, H., et al. (eds.), Moral Sentiments and Material Interests, The Foundations of Cooperation in Economic Life, Ch. 3, pp. 75–113 (2005)
30. Broom, M., Ruxton, G.D.: Evolutionarily stable kleptoparasitism: consequences of different prey types. Behav. Ecol. **14**(1), 23–33 (2003)
31. Frederickson, M.E.: Rethinking mutualism stability: cheaters and the evolution of sanctions. Q. Rev. Biol. **88**(4), 269–295 (2013)

# Towards Eusociality Using an Inverse Agent Based Model

**John C. Stevenson**

**Abstract** The emergence of eusocial species is both very rare in evolutionary history and results in remarkably successful species. Using a linear genetically-programmed agent-based model, agent rules are discovered that display behaviors characteristic of eusocial species. By holding the agents' genome constant across the colony and allowing the agents' rules to evolve, the individual behaviors exhibit phenotypic plasticity in response to environmental cues. The phenotypically driven reduction of intrinsic growth rates and the emergence of non-reproducing phenotypes both demonstrate selection pressure at the colony (system) level. Various other emergent eusocial behaviors are identified and discussed. A path forward to more capable eusocial populations and inter-colony evolution is outlined.

## 1 Introduction

Eusocial species represent a very small fraction of the total species on earth and yet they rank among the most ecologically dominant land animals by population and biomass [1]. The limited number of species that independently evolve eusociality in diverse taxa suggest this occurrence is a phylogenetically rare event and is considered "one of the great mysteries of biology". The definition of eusociality has changed since its first use in 1966 for nesting bees [3, 4]; through Wilson's classification as colonies with overlapping generations, division into reproductive and non-reproductive castes, and cooperative care for the young [5–7]; to an explicit definition that tries to incorporate the many eusocial communities in both arthropods and vertebrates [8]. For the purposes of this paper, Wilson's classification is unambiguous. Additional eusocial characteristics often found include nesting, environmental effects on reproduction rates, coexistence of different phenotypes, haplodiploiy or

J. C. Stevenson (✉)
Long Beach Institute, Long Beach, NY 11561, USA
e-mail: jcs@alumni.caltech.edu

similar reproductive strategies, and other cooperative behaviors such as group foraging and defense [9]. For those colonies whose reproductive caste is singly mated queens, all the female members of these colonies have very similar genomes; and the diverse physical and behavioral female phenotypes found within the colony are due to responses to each individual's environment (phenotypic plasticity).

Agent based models, as used in this research, are inherently social. Agents interact with and affect not only the environment but also compete and cooperate with each other. Classification of biological, sociological, and ecological models include minimal models for systems and synthetic models of systems [10]. Synthetic models of systems match the macroscopic results of the model to empirical data [11–15] and provide explanatory rules [11, 12, 16, 17]. The rules are either manually crafted or automatically generated with evolutionary algorithms [18, 19] such as those used in inverse Generative Social Science (iGSS) [20–22]. In contrast to these synthetic models, a minimal model of a system does not attempt to calibrate to an empirical objective function. Rather, a population of agents freely evolves within an environment under evolutionary selection. Some models in this category apply selection pressure exogenously [23–25]. Others apply the selection pressure endogenously within the simulation as a "struggle for existence", where more fit individuals reproduce and replace the less fit [27–29]. When applying evolutionary optimization methodology to these endogenously optimized minimal models, much of the complex algorithmic machinery used for optimizing candidate populations outside of the simulation is not required.

Within this group of minimal models of systems that evolve rules, some qualify as complex adaptive systems (CAS) which may optimize either on the level of individuals (CAS2) or as a system (CAS1)[30]. CAS2 systems often evolve into a "tragedy of the commons" thus stimulating research on cooperation. Game theory is one productive area for this research [24, 25, 32], but these games still optimize at the individual level (CAS2). True system level optimization requires individuals to reduce their own survival and reproductive success for the benefit of their community [2, 28, 30]. This research uses a minimal model for a system with endogenous evolution of genetically-programmed agent rules. The emergence of phenotypically driven reductions in intrinsic growth rates and of non-reproductive castes suggests that optimization is occurring at the colony level (CAS1). In this genetically programmed approach to the agents' rules, the agents' genome is held constant. Competition between colonies of different genomes would drive evolution of the queens' genomes, though that is beyond the scope of this paper.

The principal results of this research are the emergence of phenotypic plasticity within a colony of agents with identical and fixed genotypes [33, 34] and resultant colony-level optimizations (CAS1). Phenotypic behaviors reduced the intrinsic growth rate of the colony through competitive exclusion (benefiting both the individual and the colony) and through generation of viable populations of nonreproducing phenotypes (sacrificing individual reproductive success). A number of ancillary eusocial behaviors also emerged including stable coexistence of different phenotypes, cooperative foraging by different phenotypes, overlapping generations, phenotypic driven changes in reproduction rates, haploid reproduction, and pheno-

types that only breed and do not forage. The only defining eusocial behaviors not observed were cooperative care of the young, and nesting and its defense [2, 5–8].

The paper proceeds by first describing the underlying agent based model and its population dynamics with hard-wired agent rules and a genome that contains the relevant agent characteristics. The language of genetically programmed rules for the agents is defined, and hand-crafted programs that replicate the hard-wired rules are presented. Random initial agent programs are then allowed to evolve across the various (constant) genetic and computational capacity parameters. These results are discussed in the context of eusociality. Future research for evolving true eusocial colonies is outlined.

## 2  Models and Methods

### 2.1  Underlying Agent Based Model with Genetic Characteristics

The underlying spatial-temporal agent-based model (uABM) is based on a minimum model of a foraging system [27]. The agent characteristics that are part of the evolutionary process are defined as genes on a single chromosome which reproduces with occasional mutation (haploid parthenogenesis). These characteristics are stochastic infertility, puberty, sunk birth costs (rather than endowments), and introvert/extrovert preference. The remaining agent characteristics and landscape properties are fixed for each run. The agents interact on an equal opportunity (flat) landscape of resources. Detailed descriptions of the ABM parameters and processes sufficient to reproduce the uABM are provided here [35].

The dynamics that emerge from this simple underlying model have been shown to agree with time delayed logistic growth models for single species [35–37], stochastic gene diffusion models [35, 38], and modern coexistence theory [39, 40]. When an initial population of agents with random heterogeneous alleles is run with mutation and subjected to endogenous selection pressures of survival, the population evolves to one that is dominated by minimum infertility, minimum non-zero puberty, minimum birth cost, and introversion. These alleles represent selection at the individual level (CAS2) towards the maximum intrinsic growth rate possible. The zero puberty allele is not dominate due to spatial effects of immediate births, and introversion is preferred to avoid local resource competition. The resultant population dynamic is a tragedy of the commons [31], where the population has almost no resource reserves, mean agent lifetimes are brutally short, and extinctions are common due to environmental degradation, lack of resource reserves, and chaotic population level trajectories [41].[1]

---

[1] Discrete logistic growth equations generate population level trajectories that range from stable through oscillating and into chaotic regimes based on increasing values of intrinsic growth [36, 37, 42].

**Table 1**  Architecture and instruction set for agent programming langugage

| Name | Address | Function | | Values | Description |
|------|---------|----------|---|--------|-------------|
| nextI | 1–2 | Register | | 05–32 | Address of next instruction |
| bDir | 3 | Register | | UDLRZ | Best seen direction (Z = no data) |
| bDis | 4 | Register | | 0–9 | Best seen distance |
| bRes | 5 | Register | | 0–9 | Best seen resources |
| inst | 6–32 | Program | | UDLRMX | Executeable instruction |
| Instr | Description | Action/test | | Result | |
| U | Look up | Find cell max resource above > bRes | | Store in bDir, bDis, bRes | |
| D | Look down | Find cell max resource below > bRes | | Store in bDir, bDis, bRes | |
| L | Look left | Find cell max resource left > bRes | | Store in bDir, bDis, bRes | |
| R | Look right | Find cell max resource right > bRes | | Store in bDir, bDis, bRes | |
| M | Move | Fetch bDis, bDir, if 'Z' random values | | Move bDis, bDir | |
| X | Reproduce | Space, birth costs allow reproduction | | Place new agent in cell | |

The uABM provides the structure upon which genetic programming of the agents' behaviors is implemented. This approach presents a very large solution space of various combinations of infertility, birth cost, introvert/extrovert, and puberty alleles. Based on the cited results with the uABM using hard-wired agent rules and genetically evolving agent characteristics, the genome parameter space is reduced to only infertility and birth cost alleles. Puberty is held constant at one generation and the introvert/extrovert preference is disabled. Computation capacity of the agents adds a third parameter to the space. Haploid reproduction as clones was selected for simplicity (as exemplified by eusocial ant species Mycocepururs Smithii of Hymenoptera:Formicidae [43, 44]).

## 2.2  Agent Programming Language and Grammar

A simple language replicating the uABM agent rules was designed and integrated into an inverse ABM (iABM). Each agent has a 32 character string which contains the registers and instructions which the simulation executes on each agent's action cycle. Five characters are used for registers leaving up to 27 characters for the program. These instructions are described in Table 1. The action cycle for this iABM is depicted as a flow chart in Fig. 1.

**Fig. 1** Action cycle for iABM

The number of instructions that can be executed per each agent's action cycle, called computation capacity, is part of the parameter space that is surveyed. Foraging gains, metabolic costs, and deaths occur during the move instruction. Since multiple moves may occur during one action cycle, each move instruction triggers foraging at the new location and incurs the metabolic cost. Birth decisions and associated costs occur during the reproduction instruction. If the action cycle ends without at least one metabolic resource cost, one is applied. With this genetic programming approach, the hard-wired rules of the uABM can be replicated with a computation capacity of six steps. These programs ("classic" phenotypes) contain 6 instructions in one action cycle that look in four directions, move, and reproduce (e.g. UDLRMX and 23 other versions of look ordering). For hard-wired rules, a tie for the best direction and distance is broken randomly. For the genetic programming version (iABM), the

first instruction is equally seeded with each of the four directions. Other than for these ties, the order of look instructions (before a move) is not functionally significant. If a move is targeted to a location that is no longer valid (occupied either due to a random move based on no look data collected since the last move, or from outdated look data from a previous action cycle) the agent does not move. The results, with initially seeded classic phenotypes for simulations spanning the genome alleles, has indistinguishable population dynamics and agent metrics from the uABM. These classic phenotypes often emerge as good solutions and, surprisingly, are sometimes competitively excluded, usually by pairs of cooperating phenotypes.

## 2.3   Methods

All runs are initiated with a population of 400 programs with random instructions of random length. Different seeds generate different initial phenotype populations and resultant population trajectories. Some genetic allele combinations are so challenging that only a few of the initial random sets of phenotypes are able to generate viable, reproducing populations. Figure 2a gives the fraction of the initial random population that survives through the initial population minimum and is fertile. Each point is the mean fraction of viable surviving phenotypes over 40 differently seeded runs with 400 initial random phenotypes each. Simulation runs are generally stopped at either 10,000 or 50,000 generations, orders of magnitude past the attainment of steady population levels. Phenotype evolution occurs continuously through out the simulations so the stopping point is somewhat arbitrary. Events of interest or long running trends receive longer run times (Figs. 2b and 3).

The optimization of the instructions defining the agents' behavior through genetic programming is straightforward. Genetic algorithms and genetic programming [19, 45] have a large body of techniques for shaping evolving populations [46–50] and for multi-objective optimization [51]. Since the genetic programs that form phenotypic behavior are selected and propagated continuously throughout the simulations based on a "struggle for existence", the complex art of exogenous population optimization is avoided. As Fig. 2a shows, random initial instruction sets over tens of runs are sufficient to generate viable and interesting phenotypes. When an agent reproduces, a single point mutation will occur in the daughter agent at a constant probability $\mu$ per reproduction. If a mutation occurs, a location in the program and a type of mutation are chosen randomly. Three mutation types are implemented: flip to a different random instruction; insert a new random instruction if memory space allows; or knockout the instruction (if the program is longer than one instruction).

**Fig. 2** Frequency of occurrence of phenotypes of interest. **a** The percent fraction of initial, random phenotypes that survive through the initial population minimum. **b** The percent fraction of non-reproducing phenotypes at the final generation for all infertility and birth costs

## 3 Eusocial Behaviors

The emergence of phenotypic plasticity displayed a surprising number of behaviors characteristic of eusocial communities. The consistent emergence of viable, non-reproductive phenotypes is a significant milestone for eusocial behavior. Phenotypically driven changes in growth rate modify the rate set by the colony's genome and result in both higher intrinsic growth rates benefiting the individuals (CAS2); and lower intrinsic growth rates benefiting both individuals and colony (CAS1). The ability of different phenotypes to competitively coexist in accordance with modern coexistence theory [39] enables most eusocial behaviors. Coexisting phenotypes support caste emergence, cooperative foraging, higher colony fitness, stable sub-populations of sterile phenotypes, and influence population volatility.

## 3.1 Populations with Significant Fractions of Non-reproducing Phenotypes

Division of reproductive labor is one of the defining characteristics of an eusocial society. Figure 2b presents the statistics on the fractions of the final populations which are non-reproducing but viable. These statistics are taken over the constant alleles of infertility and birth cost and presented by computation capacity. Figure 3 exemplifies phenotype population trajectories for two representative simulations with differing birth costs where the non-reproducing fraction of the population was greater than half and rose over time. The commonality and viability of these phenotypes suggest

**Fig. 3** Division of reproductive labor. **a** The emergence of a large population of non-reproducing phenotypes (no X) for infertility and birth cost 1, mutation rate 0.01, and computation capacity 2. **b** The emergence of a large population of non-reproducing phenotypes (no X) for infertility 1 and birth cost 5, mutation rate 0.01, and computational capacity 2

two important points. One, there is selective pressure for the emergence of non-reproducing phenotypes which demonstrates system level optimization (CAS1) since these phenotypes never reproduce. These phenotypes consume resources and space without reproducing, which can benefit the colony by reducing the overall intrinsic growth rate, helping to avoid oscillatory and chaotic population trajectories. Though non-reproducing castes are always justified with "cooperative care of the young", these experiments suggest that there are other benefits to the colony [2].

## 3.2 Phenotypically Driven Colony Growth Rates

Phenotypic behavior can increase or decrease the colony's intrinsic growth rate from that specified by the genome. The uABM replicates discrete logistic growth (with time delay) which has transitions to oscillating and chaotic population level regimes with increasing intrinsic growth rates driven by the allele values. Natural selection at the individual level under these conditions drives toward higher intrinsic growth rates (CAS2). But with a constant colony genome of intrinsic growth, phenotypic behavior often decreases the intrinsic growth rate pushing the colony into more stable regimes.[2] By pushing the colony population dynamics into more stable regimes, the colony benefits by avoiding a tragedy of the commons, chaotic exclusion of more

---

[2] During the initial growth phase into a rich landscape with few agents, phenotypes are selected for greater intrinsic growth.

fit phenotypes, and potential extinction (CAS1). Specific examples are discussed in detail in the following section. Adaption of the intrinsic growth rate of a colony to environmental conditions through phenotype plasticity is a characteristic of eusociality.

## 3.3   Coexistence and Competitive Exclusion

In many allele and seed configurations, two or more competing phenotypes will coexist, generating colony fitness that neither would be capable of alone. Other times, a well established resident will be excluded by an invading new mutation. Figure 4 presents examples of both. These two examples also provide an excellent demonstration of the wide variety of solutions that will emerge based solely on different seeding of random sets of initial instructions. In both cases, the early resident phenotypes have high intrinsic growth rates and generate high population level volatility which, by pushing the dynamics into chaotic regimes, are less fit once the landscape's carry capacity is reached and are eventually excluded. The population in Fig. 4a with only one resident phenotype has two clear exclusion events where a new mutant invades and quickly excludes a resident population [39]. These new mutants are both single instruction flips. The first exclusion event occurs around generation 1500 when a reproduction instruction (X) at position 20 mutates to a move (M) which pushes the population dynamics out of a chaotic regime. The reduction in reproduction rate with an increase in foraging out competes its parent. The second exclusion event around generation 7700 is a single mutation at position 16 from a right look (R) to a left look (L). This mutation changed the ratio of left looks to right looks in the phenotype from $\frac{4}{2}$ to $\frac{5}{1}$ producing a relatively more fit left sweeper. In Fig. 4b the early resident phenotype is again generating chaotic population dynamics and is again quickly excluded, this time by a pair of phenotypes working together to generate a population that has comparable mean but lower volatility. The first pair of coexisting phenotypes, appearing around generation 2000, sweeps east-north-east (RMXURM) and broadly south with opportunistic moves to east or west (DMXLRM). At around generation 8200 a single flip mutation of the broadly south sweeper at position 5 from a right look (R) to a down look (D) proves more fit than the east-north-east sweeper cooperating with the other phenotype. The new pair, parent and child, sweeps both broadly south and south-south-west. Both these paired sweeping patterns are suggestive of cooperative foraging.

## 4   Discussion and Future Work

The discovery of genetically programmed agent behaviors in a spatial-temporal agent-based minimal model of a system has demonstrated the emergence of creative and novel agent behavior rules relevant to eusocial societies. Phenotypic plasticity, for one, opens the door for both eusociality and inter-colony evolution. The

**Competitive Coexistence and Exclusion**



**Fig. 4** Competitive coexistence and exclusion two differently seeded solutions that emerge for a constant genome with infertility and birth cost 1, and computational capacity 3. **a** The single mutation of a right look to a left look after generation 7000 drives this invader to exclude the previous resident phenotype. **b** Exemplifies coexistence between a phenotype sweeping broadly south with, first, one that looks and moves ENE, and then replaced by a SSW sweeper

emergence and viability of non-reproducing phenotypes, a necessary but not sufficient defining behavior for eusociality, suggests selection pressure at the colony level (CAS1). Other phenotypically-driven reductions in intrinsic growth rates benefit both the individual (CAS2) and the colony (CAS1). Social insects Many behaviors characteristic of eusocial societies are shown to have emerged from random initial populations of programs whose agents all posses the same colony genome. Cooperation through coexistence leads to higher colony fitness. Non-reproducing phenotypes emerged and increased to majority representation in many colonies. Phenotypic plasticity significantly changed the intrinsic growth rate of the colony, moving it into or out of oscillatory or chaotic population regimes. These changes in intrinsic growth rate were often achieved by cooperating phenotypes.

Numerous examples of distinct, novel and informative agent behaviors based on environmental conditions exhibited phenotypic plasticity. Classic phenotypes often emerged for computational capacity 6 but were often competitively excluded by pairs of cooperating phenotypes.

Emergence of conventionally defined eusocial colonies using this model will require the addition of local sharing of resources (cooperative care of young) and sensing local neighbors' colony genome (friend/foe) which, when combined with exploitation of the introvert/extrovert gene, may generate nesting behaviors (philopatry). The current structure of this iABM with a separate queen's genome for each colony coupled with phenotypic plasticity through evolving agent rules supports inter-colony competition and evolution of the colonies' genomes.

# References

1. Wilson, E.O.: Social insects. Science **172**(3981), 406 (1971)
2. Howard, K.J., Thorne, B.L.: Eusocial evolution in termites and hymenoptera. Biology of Termites: A Modern Synthesis, pp. 97–132 (2011)
3. Batra, S.W.: Nests and social behavior of halictine bees of India (hymenoptera: Halictidae). Indian J. Entomol. **28**, 375 (1966)
4. Michener, C.D.: Comparative social behavior of bees. Ann. Rev. Entomol. **14**(1), 299–342 (1969)
5. Wilson, et al., E.O.: The Insect Societies. Harvard University Press, Cambridge, Massachusetts, USA (1971)
6. Wilson, E.O., Hölldobler, B.: Eusociality: origin and consequences. Proc. Nat. Acad. Sci. **102**(38), 13367–13371 (2005)
7. Ward, P.S.: The phylogeny and evolution of ants. Annual review of ecology. Evol. Systematics **45**, 23–43 (2014)
8. Crespi, B.J., Yanega, D.: The de nition of eusociality. Behavioral Ecol. **6**(1), 109–115 (1995)
9. Friedman, D., Johnson, B., Linksvayer, T.: Distributed physiology and the molecular basis of social life in eusocial insects. Hormones and behavior **122**, 104757 (2020)
10. Roughgarden, J., Bergmen, A., Hafir, S., Taylor, C.: Adaptive computation in ecology and evolution: a guide for future research. Adaptive individuals in evolving populations. SFI Studies in the sciences of complexity **26** (1996)
11. Patel, A., Crooks, A., Koizumi, N.: Spatial agent-based modeling to explore slum formation dynamics in ahmedabad, India. GeoComputational Analysis and Modeling of Regional Systems, pp. 121–141 (2018)
12. Grimm, V., Revilla, E., Berger, U., Jeltsch, F., Mooij, W.M., Railsback, S.F., Thulke, H.H., Weiner, J., Wiegand, T., DeAngelis, D.L.: Pattern-oriented modeling of agent-based complex systems: lessons from ecology. Science **310**(5750), 987–991 (2005)
13. Wiegand, T., Jeltsch, F., Hanski, I., Grimm, V.: Using pattern-oriented modeling for revealing hidden information: a key for reconciling ecological theory and application. Oikos **100**(2), 209–222 (2003)
14. Bianchi, C., Cirillo, P., Gallegati, M., Vagliasindi, P.A.: Validating and calibrating agent-based models: a case study. Comput. Econ. **30**, 245–264 (2007)
15. Heckbert, S., Baynes, T., Reeson, A.: Agent-based modeling in ecological economics. Ann. N. Y. Acad. Sci. **1185**(1), 39–53 (2010)
16. Brugnera, M.D.P.E., Fischer, R., Taubert, F., Huth, A., Verbeeck, H.: Lianas in silico, ecological insights from a model of structural parasitism. Ecol. Modell. **431**, 109159 (2020)
17. Epstein, J.M.: Agent-based computational models and generative social science. Complexity **4**(5), 41–60 (1999)
18. Holland, J.H.: Adaptation in Natural and Artificial Systems: An Introductory Analysis with Applications to Biology, Control, and Artificial Intelligence. MIT Press (1992)
19. Koza, J.: On the programming of computers by means of natural selection. Genetic programming (1992)
20. Vu, T.M., Probst, C., Epstein, J.M., Brennan, A., Strong, M., Purshouse, R.C.: Toward inverse generative social science using multi-objective genetic programming. Proceedings of GEECO, pp. 1356–1363 (2019)
21. Gunaratne, C., Garibay, I.: Alternate social theory discovery using genetic programming: towards better understanding the articial anasazi. Proceedings of the Genetic and Evolutionary Computation Conference, pp. 115–122 (2017)
22. Greig, R., Arranz, J.: Generating agent based models from scratch with genetic programming. ALIFE 2021: The 2021 Conference on Artificial Life. MIT Press (2021)
23. Fogel, D.B.: Evolving behaviors in the iterated prisoner's dilemma. Evol. Comput. **1**(1), 77–97 (1993)
24. Lindgren, K., Nordahl, M.G.: Evolutionary dynamics of spatial games. Physica D Nonlinear Phenomena **75**(1–3), 292–309 (1994)
25. Miller, J.H.: The coevolution of automata in the repeated prisoner's dilemma. J. Econ. Behavior Org. **29**(1), 87–112 (1996)

26. Gause, G.F.: The Struggle for Existence. Williams and Wilkins (1934)
27. Epstein, J.M., Axtell, R.: Growing Artificial Societies from the Bottom Up. MIT Press (1996)
28. Pepper, J.W., Smuts, B.B.: The evolution of cooperation in an ecological context: an agent-based model. Dynamics in Human and Primate Societies: Agent-Based Modeling of Social and Spatial Processes, pp. 45–76 (2000)
29. Epstein, J.M.: Zones of cooperation in demographic prisoner's dilemma. Complexity **4**(2), 36–48 (1998)
30. Wilson, D.S., Kirman, A.: Complexity and Evolution: Toward a New Synthesis for Economics, vol. 19. MIT Press (2016)
31. Ostrom, E.: Governing the Commons: The Evolution of Institutions for Collective Action. Cambridge University Press (1990)
32. Axelrod, R., Hamilton, W.D.: The evolution of cooperation. Science **211**(4489), 1390–1396 (1981)
33. West-Eberhard, M.J.: Phenotypic plasticity and the origins of diversity. Ann. Rev. Ecol. Syst. **20**(1), 249–278 (1989)
34. DeWitt, T.J., Sih, A., Wilson, D.S.: Costs and limits of phenotypic plasticity. Trends Ecol. Evol. **13**(2), 77–81 (1998)
35. Stevenson, J.C.: Agentization of two-population driven models of mathematical biology. Proceedings of the 2021 International Conference of the CSSSA (2021)
36. Murray, J.D.: Mathematical Biology. Springer (2002)
37. Kot, M.: Elements of Mathematical Ecology. Cambridge University Press (2001). https://doi.org/10.1017/CBO9780511608520
38. Ewens, W.J.: Mathematical Population Genetics: Theoretical Introduction, vol. 1. Springer (2004)
39. Chesson, P.: Mechanisms of maintenance of species diversity. Ann. Rev. Ecol. Syst. **31**(1), 343–366 (2000)
40. Stevenson, J.C.: Competitive exclusion in an articial foraging ecosystem. arXiv:2203.02814 (2022)
41. Stevenson, J.C.: Dynamics of wealth inequality in simple articial societies. Advances in Social Simulation, pp. 161–172. Springer (2022)
42. Liz, E.: Delayed logistic population models revisited. Publicacions matematiques, pp. 309–331 (2014)
43. Rabeling, C., Gonzales, O., Schultz, T.R., Bacci, M., Jr., Garcia, M.V., Verhaagh, M., Ishak, H.D., Mueller, U.G.: Cryptic sexual populations account for genetic diversity and ecological success in a widely distributed, asexual fungus-growing ant. Proc. Nat. Acad. Sci. **108**(30), 12366–12371 (2011)
44. Himler, A.G., Caldera, E.J., Baer, B.C., Fernandez-Marin, H., Mueller, U.G.: No sex in fungus-farming ants or their crops. Proc. R. Soc. B Biol. Sci. **276**(1667), 2611–2616 (2009)
45. Holland, J.H.: Building blocks, cohort genetic algorithms, and hyperplane-defined functions. Evol. Comput. **8**(4), 373–391 (2000)
46. Gupta, D., Ghafir, S.: An overview of methods maintaining diversity in genetic algorithms. Int. J. Emerg. Technol. Adv. Eng. **2**(5), 56–60 (2012)
47. Poli, R., McPhee, N.F., Vanneschi, L.: Elitism reduces bloat in genetic programming. Proceedings of the 10th Annual Conference on Genetic and Evolutionary Computation, pp. 1343–1344 (2008)
48. Esparcia-Alcázar, A., Sharman, K.: Phenotype plasticity in genetic programming: A comparison of darwinian and lamarckian inheritance schemes. Genetic Programming: Second European Workshop, EuroGP'99 Göteborg, Sweden, May 26–27, 1999 Proceedings 2, pp. 49–64. Springer (1999)
49. La Cava, W., Helmuth, T., Spector, L., Danai, K.: Genetic programming with epigenetic local search. Proceedings of the 2015 Annual Conference on Genetic and Evolutionary Computation, pp. 1055–1062 (2015)
50. Diaz-Gomez, P.A., Hougen, D.F.: Empirical study: Initial population diversity and genetic algorithm performance. Artif. Intell. Pattern Recogn. **2007**, 334–341 (2007)
51. Rodriguez-Vazquez, K., Fonseca, C.M., Fleming, P.J.: 'identifying the structure of nonlinear dynamic systems using multiobjective genetic programming. Man, and cybernetics—part A: systems and humans. IEEE Trans. Syst. **34**(4), 531–545 (2004)

# Utilizing the Full Potential of Norms for the Agent's Decision Process

**Christian Kammler** ⓘ**, Frank Dignum** ⓘ**, and Nanda Wijermans** ⓘ

**Abstract** Norms are a crucial part of human behavior that received a lot of attention within the social simulation community. However, some aspects—up until now—have not been addressed in existing agent architectures, such as their motivational aspects and their importance and impact in planning and action selection. In this paper we present an agent architecture capable of grasping this potential of norms. We use perspectives to reflect how different people engage with a norm, and how it effects their long-term goals, their planning, and course of action. Our architecture is capable of having fast habitual-like behavior, as well as more complex deliberation if necessary.

**Keywords** Norms · Values · Social rules · Social simulation

## 1 Introduction

Norms are a crucial part of human behavior, and influence it in a variety of ways and multiple levels [9]. While norms have received a lot of attention within the social simulation community, see e.g. [1, 3, 9, 11, 15, 21], their motivational aspects as well as their importance and full impact in planning and action selection have not been incorporated in existing agent architectures. Such realistic—human like—behavior is especially important when simulating policies (legal norms) and the reactions to them, as the discussion around COVID-19 showed [11, 18].

C. Kammler (✉) · F. Dignum
Department of Computing Science, Umeå University, Umeå, Sweden
e-mail: christian.kammler@umu.se

F. Dignum
e-mail: frank.dignum@umu.se

N. Wijermans
Stockholm Resilience Centre, Stockholm University, Stockholm, Sweden
e-mail: nanda.wijermans@su.se

To achieve norm realism in models, agents are required to not only see if the consequences of following or breaking the norm are desirable for them, but also how they are motivated [5, 21] to circumvent the norm. Furthermore, a new norm often interacts with existing normative structures as well (such as socially accepted behavior—social norms), which might cause that to change. Since norm breaking is an important part of norm change [2, Chap. 5], the agents will also need a flexible way to deal with norm violations, where they can chose at multiple stages in their decision making whether or not to violate a norm as well as to react to other agent's norm violations. As people react differently to norms, focusing on the parts that are relevant for them, we also need to include different perspectives on norms [17]. Each perspective reflects their (priority in) own goals and values, and available actions [17].

To combine these requirements in an agent (i.e. motivation of norms, flexible way to deal with norm changes and violations, and perspectives on norms) we connect norms within the different stages in the agent's decision making process: (a) goal selection, (b) planning the course of action, and (c) action selection in a changing environment. While approaches are existing for each of these different parts, e.g. [11] for goal and action selection, or [22] for planning, they have drawbacks which make them not suitable for our approach, such as an only implicit representation of norms and no planning capabilities [11], or not taking the motivational aspects of norms into account [22], which we both require.

In this paper we propose a norm and value based agent deliberation process (see Sect. 2) and operationalise this in a novel agent architecture (see Sect. 3). Creating agents that are able to reflect the influences and effects norms can have on one's deliberation processes and consequently make situated decisions. We also enable agents to have different levels of complexity in their deliberation based on the availability of actions and plans. This involves quick and habitual behavior given a familiar situation, as well as to bring in more complex deliberation when necessary, such as finding alternative actions or planning for a goal. To demonstrate our deliberation process works, we use an example of a restaurant size based restriction on the number of guests [17, 18], see Sect. 4. We conclude with a discussion in Sect. 5.

## 2 Our Approach to Normative Human Decision Making in the Context of Norms

To develop a norm and value based agent deliberation process, we first lay the basis on how different people engage with a norm and depart from there with two main building blocks: perspectives on norms, and norms.

## 2.1 Perspectives on Norms

We use perspectives on norms to describe individual differences in norm engagement. Norm engagement concerns a focus on only the relevant parts of the norm and consequently used differently in our decision-making. For example: Restaurant owners focus on the financial impact of a size-based restriction norm, whereas guests focus on the social impact of that same norm [17]. The concept of perspectives describes the individual differences defined as: "*A perspective is specified by goals (G), available actions (A), effects of those actions (EoA), social affordances (SocAffs), and priorities in values (PrioV)*" [17, p.142].

Among the goals of a perspective we distinguish *maintenance goals* and *achievement goals*. *Actions* and *social affordances* of a perspective are mapped into physical and social actions. Physical actions require one or more objects ($o_i, o_j, ...$) to perform the action. Actions also have a pre-condition that needs to be met so the action can be executed. Actions can also change the physical state of the simulation, which is called the result of the action [17]. Social actions are the social effects of the physical actions, i.e. the social affordances [17]. For example: Sitting together with friends at a table (physical action) has in the social dimension the effect (social effect) of socializing with one's friends (social action), i.e. it socially affords the action.

Norms influence actions in the way that some actions might become obligatory or forbidden. Firstly by identifying the actions that are affected by a norm are done using the object of the norm ($I_{Object}$); by checking if $o_i == I_{Object}$ holds, where $o_i$ is the object required by the action; and finally, actions pro-/demote values (being more or less important in a given situation).

## 2.2 Maintenance and Achievement Goals

Approaches based on motivation theories have been shown to be successful to model and explain reactions to policies, see e.g. [8] (smoking ban) and [11] (COVID-19). We adjust the watertank model (tanks run low and can be filled up based on the satisfaction of the drives) originated by [12] and successfully used in [11, 16]. Because different people have different drives, we define different drives for different perspectives. We separate these drives into maintenance goals (states the agent wants to maintain, watertanks) and achievement goals (states the agent wants to achieve to fill up the watertanks).

**Maintenance goals** reflect the states that the agents want to maintain [8, 13], e.g., always keep the restaurant room temperature in a certain range. Every perspective has their own maintenance goals. However, all people want to have a functioning society which we represent as the maintenance goal of conformity. For a maintenance goal to be salient we use a threshold between zero and one, where the height of the threshold reflects the importance of the particular maintenance goal [16]. The higher the difference between fill level and threshold, the greater the urgency of that

maintenance goal [16]. Over time satisfaction of the goals decreases with a constant amount [11, 16]. Furthermore, the higher importance of the goal, the higher the constant depletion to further reflect the importance of a certain watertank.

**Achievement goals** are the concrete goals, such as reserving a table, the agent wants to achieve [8, 13]. Achieving these goals satisfy maintenance goals and pro-/demote values (see Sect. 2.3). These goals are achieved by (sequences) of actions, e.g., reserving a table is reached by the action to call the restaurant. Note that it is possible that the goal is pro-/demoting multiple values. Connecting values with actions (through executing goals) enables the agent to reason about which goals are aligning with their values.

## 2.3   *Values*

Values provide an evaluation mechanism for actions and events [11, Chap. 2]. They motivate us to reach states which are aligning with our values [11, Chap. 2]. Values can be seen as a form of ordering of preferences. Actions and world states that promote the person's values are preferred over other. states [8].

To reflect values, we use the Schwartz value system [23]. We adopt the ten universally recognized values: self-direction, stimulation, hedonism, achievement, power, security, conformity. tradition, benevolence, and universalism. It has successfully been implemented in several agent-based social simulations to model and explain reactions to policies, see e.g. [7, 11], and for its establishment [16].

We use values to determine which achievement goals are important, which actions are most desirable to take, and which maintenance goals are more urgent to satisfy, given their influence on the threshold and depletion. To do this, we adopt Heidari's [16] work who connected Schwartz values [23] to concrete actions. Each perspective has its own priority of values (PrioV) [17] which are constant. These are then compared with actions and goals to see which pro-/demote their desired values.

The priority of values (PrioV) are defined as: a total set of values $V = \{V_1, ..., V_{10}\}$ in the simulation. For each perspective they are ordered differently, whereby the order in the list determines the importance. We ensure in this ordering that opposite values cannot be of similar importance. To use them in calculations, each value is mapped to a number $n \in [0; 1]$. The higher the importance, the higher the value of n.

Restaurant owners with a PrioV in power and achievement, for example, will be strongly money driven (i.e., the maintenance goal of `making money` is very urgent for them to satisfy), and therefore want to achieve goals and take actions, such as increasing the price or reducing the costs of their meals, to make as much money as possible [17]. Since we are also going to attach values to norms, agents can use their priority in values to decide whether to violate a norm or not.

To actually make a decision on their next desired action, the agent looks first if it is forbidden by a norm. If the norm promotes the same values which are important for the agent, it will comply to the norm. If it is not the case, the agent compares

the cost of complying to the norm (the loss in satisfaction gained) vs the cost of not complying to the norm (cost on conformity). The selection is then based on an `argmax function` with the PrioV acting as multipliers.

## 2.4 Norms

To model behavior and responses to norms, especially policies, it is necessary to take into account the motivational parts of norms [5, 21]. People will try to find a way around the (legal) norm, thereby also generating new behavior [19].

To formalize norms, we use the ADICDlRO framework [17] (an extension of the ADICO grammar by [6] allowing agents to reason about norms based on their perspective. Norms in the ADICDlRO framework are defined by [17, 18]: A specifies the agent group which is responsible to adhere to the norm. D is the deontic part of the norm, and together with the aIm (I), split into action ($I_{Verb}$) which the norm is targeting and the object ($I_{Object}$) of the norm, they form the `{fulfilment, violation}` condition of the norm. C defines the contexts in which the norm is active and not active, therefore representing the `{activation, deactivation}` condition of the norm. This is different from the deadline element (Dl) which states when the norm is supposed to be fulfilled. If a norm is supposed to be always fulfilled then the deadline is set to time 0. The repair part (R) of the framework defines the action(s) to 'undo' the potential breaking of the norm, and the 'Or else' (O) specifies the punishment of the norm violation.

An important addition of the ADICDlRO framework is that norms pro-/demote values, based on the purpose they fulfill [18], allowing agents to reason if the norm is important to them, based on their PrioV. For example, the size-based restriction on the allowed number of guests in the restaurant is introduced to combat the spread of COVID-19, and thus promotes safety. However, it also has a potential negative impact on the revenue by the restaurant owner and thus, demotes power and achievement. This added value dimension makes it now possible for agents to deliberate if the norm is important for them or if they want to violate it, based on how it aligns with their values.

A complete instantiation of our Example norm—the size-based restriction on the number of guests in a restaurant—can be as follows, assuming it is active all the time (deactivation condition == none), and available at the start of the simulation (timestep == 0) (based on [17, 18]): A restaurant owner (A) must not (D) have ($I_{Verb}$) more than X guests in their restaurant ($I_{Object}$, D + I = $\{\#guests \leq X, \#guests > X\}$) at all times (`Condition = {timestep == 0, none}`), effective immediately (D, `timestep==0`), or else they will be fined with 5000\$ (O). This norm promotes the value of safety (`promoted values`), demotes the value of power (`demoted values`), and any violation of this norm can be undone by making the guests that are too much leave (R).

## 2.5   *Planning with Values*

While individual actions give us options to choose from, they will not always be able to achieve the agents goals in one step. This requires that we have some mechanism for linking multiple actions in order to achieve a goal. The agent is given planning capabilities by plan patterns [10]. Plan patterns have the following advantages: They help to find a balance between pro-active goal directed behavior and reactive situational behavior. Agents also do not need to have fully fleshed out plan, can plan from landmark to landmarks which can be discussed in conversations with stakeholders, enabling participatory modeling.

Formally, plan patterns describe sets of sequences of actions, defined in the terms of the landmarks. Landmarks are fixed points within a plan that must be achieved along the way [18], e.g., paying when visiting a restaurant. To select a plan pattern for the current achievement goal, the agent looks at the last landmark in the sequence to see if it is the goal state. Alternatively, we can explicitly label the plan pattern with the goal it achieves. Now the agent only has to find a plan for the next landmark, which means that the landmark basically acts as a sub-goal for the current achievement goal. For example: if the next landmark is to pay at the restaurant, the agent just has to find a plan to achieve the "has paid" state. This can either be in form of concrete actions or another plan pattern which is then iterated over until only concrete action remain.

To make plans to achieve goals and landmarks, we use goal-oriented action planning (GOAP) [20] which is based on STRIPS (Stanford Research Institute Problem Solver) [14]. To plan, the agent looks at the goal state (e.g. having payed), and creates a state trajectory backwards to the current state. Actions are thereby used to transition between the states. If no state trajectory can be created, the goal or landmark is not achievable. In this planning, the agent also decides on violating or adhering to norms.

## 3   The Deliberation

Our agent deliberation architecture integrates the elements discussed in the previous section. It consists of four levels: the base (internal agent drives and norms), simple, medium and high complex deliberation. Spanning a habitual form of deliberation (simple) to handling situations where alternative actions need to be found (medium) to also involve planning or new goals (complex).

## 3.1   *The Base Level*

Figure 1 reflects the base for the agent's deliberation process described as follows.

1. **Perspectives**: Each perspective is attached to a set of agents, (bottom rectangle of the Fig. 1) determine their maintenance goals, achievement goals, priority in values, which actions they have available and how they look at norms.

**Fig. 1** Base level—action DB and norm DB are public, the rest is private

2. **Watertank Manager (WT) + Achievement Goal DB (blue)**: Each agent's maintenance goals get satisfied differently. To take care of this, the maintenance goal manager handles the changes in fill levels of each watertank, based on the depletion and satisfaction gain from the previous deliberation. The achievement goal database stores all available achievement goals the agent has.

3. **Priority in Values (PrioV, red)**: PrioV are used where planning, goal or action selection are involved. The numerical values used in these calculations are constant, and the higher priority the value has, the higher its numerical value is. Agents with the same perspective have the same PrioV.

4. **Norm Manager + Norm database (green)**: The norm manager (green rectangle) handles the norm (de)activation, and stores norm violations in the norm database (green). Norms are either active or not based on the current world state. (De)Activation is formalised as a set of $N$ *norms* $\{N_1, ..., N_n\}$ in the simulation. Furthermore, S is the current world state. To check if a norm has to be (de)actived, the norm manager checks for every currently (de)active if the norm's (de)activation condition is part of S. If this is the case, the norm is going to be (de)activated. To see if a norm is (de)active, we use a flag to each norm (e.g. a bool variable) signaling whether it is (de)active.

5. **Plan and Action Database (no color)**: The plan database contains the plans (sequences of actions to achieve an achievement goal) comprised out of a (sub)set of actions that exist in the action database. Each agent makes their own plans, based on the individual situations they are in. The action database holds both physical and social actions, this is situated external to the agent as action databases are connected to the different perspectives which all agents with the same perspective may make use of.

## *3.2   Simple Deliberation*

The blue background in Fig. 2 visualizes simple deliberation, reflecting "business as usual". Such habitual behavior occurs when a goal, plan, and course of action that matches the current situation for the agent are there.

The entry point is the **assessing situation** step. The situation the agent is currently in is stored in a list *currentSituation*. This list is used in the current deliberation cycle. $currentSituation = [ID, G, P, l_{phys}, l_{soc}, t_{phys}, t_{soc}, PrioV, WT - Level, N_{active}]$, with: ID = ID of the current step, G = the current goal, P the current step in the plan to achieve that goal, $l_{phys}$ = the current physical location, $l_{soc}$ = the current social location (social meaning of the physical location), $t_{phys}$ = the current physical time, $t_{soc}$ = the current social time (social meaning of the current physical time, PrioV = the agent's priority in values, WT-Level = the agent's watertank levels, and $N_{active}$ = the current active norms. The agent checks if the goal is reached. If the goal is reached, the agent has to generate a new goal. Otherwise, the agent proceeds with its current plan.

To see if an **external trigger** happened, the agent compares the *previousSituation*—residing in the simulation states database—with the *currentSituation* by detecting a change in the social time ($t_{soc}$) or the social location ($l_{soc}$) or whether a new norm is active in the *currentSituation* compared to the *previousSituation*. If this is the case, we check if the agent has been in the situation before. Otherwise, the agent prepares its next step (**What is next**).

The next step can be either an action or a landmark. If it is an action, the agent has to see if the pre-conditions hold, and whether there is any norm conflict. Each consideration results in different deliberation steps: finding a plan for a landmark (in case of a landmark), find an alternative action (in case of a norm conflict), or simply executing the action.

## *3.3   Medium Complex Deliberation*

The beige background in Fig. 2 visualizes he medium complex deliberation of the architecture. It extends the simple deliberation with the need for an agent to find an alternative action (from the action database) when it's current action (prepare the next step) is not executable anymore. The agent does this by looking for an action that has a similar post-condition as the previously desired action. If such an action is found, the agent checks for this action if the pre-conditions for the newly found action hold. The action will be executed if the action is executable, and no norm conflicts exist or violating a norm is cheaper than adhering to it. In any other case, the agents keeps querying the action database until it finds an action with a mechanism to avoid an infinite loop.

**Fig. 2** Deliberation architecture with the base level at the bottom. Simple deliberation has a blue, medium-complex a beige background, and complex a green background

## 3.4 Most Complex Deliberation

Lastly, the green background in Fig. 2 shows the most complex deliberation, involving planning and achievement goal selection. For planning we use GOAP (Goal-Oriented Action Planning, a simplified STRIPS-like planner mainly used for real-time control of autonomous agents in game development) [20] to achieve the current goal. To select an achievement goal, the agent is calculating the overall satisfaction gain gained by each achievement (not currently forbidden by a norm), and then performs an `argmax` over the satisfactions gains and selects the highest one. This is followed by creating a state trajectory (plan) from the goal state to the current state to achieve the selected goal.

## 4 Example

To show how our architecture works and compare it to other approaches, we use the following norm as an example: A size-based restriction of the concurrently allowed number of guests in a restaurant, promoting the value of safety. While multiple groups are affected by this norm, we focus on the guests and restaurant owners, each representing one perspective. The maintenance goals of the restaurant owner are to make money (WT1), and conform to the norms (WT2) with their PrioV = *power* and

*achievement*, i.e. they are strongly profit driven. The maintenance goals of the guests are to have pleasure (WT1), and to conform to the norms (WT2) with their PrioV = *benevolence*, i.e.,they are very social driven. Note, we do not strive for completeness with this example, as many aspects could be included, however to demonstrates how beneficial our approach is compared to other existing approaches, this small subset of possibilities suffices. Furthermore, we assume that the norm is not active in the beginning.

The **restaurant owner** is assessing its *currentSituation: [2 (ID), "Business as usual" (Goal), "daily check up" (Step in the plan), restaurant (phys. location), working place (soc. location), evening (phys. time), working time (soc. time), power & achievement values (PrioV, their two most important ones), (0.5, 0.8) (WT1-level, WT2-level), size-based restriction (active Norms)]*. Comparing this with the *previousList = [1, "Business as usual", open restaurant, restaurant, working place, early evening, working time, (power, achievement), (0.7, 0.85), None]*, the restaurant owner realizes that an external trigger happened (the new norm (size-based restriction) is now active). Querying the *Assessed Situations DB* yields no result, meaning that the restaurant owner has not been in this situation before. Thus a new achievement goal has to be selected. After performing an `argmax` over all non-forbidden achievement goals to find the highest overall needs satisfaction, the restaurant owner decides to lower the variable costs (the costs that they can influence), as this fills up its watertanks (in this case the maintenance goal of making money with the restaurant) the most. Next is to make a plan for this goal. The restaurant owner finds a plan pattern for this goal in the plan database. The landmarks for this plan pattern are to have all dishes with a lot of sauce identified, and the use of cheaper ingredients (as the sauce can cover the taste). To do so, the agent is looking into their action database and selects the action to filter their menu to collect all dishes with a lot of sauce. This action is not forbidden by the norm, and thus added to the plan. Also, since this is only one action, it will directly be executed. In the following step, the *currentSituation* has the following values *currentSituation = [3, Lower variable costs, use cheaper ingredients, restaurant, working place, late evening, working time, (power, achievement), (0.6, 0.82), size-based restriction]*. Since no external trigger happened (same social location and time, and no new norm active), the restaurant owner is preparing the next step in the deliberation, and then goes through the same steps as above.

The **guests** have a similar deliberation: *currentSituation:[4 (ID), "Dinner with friends" (Goal), "check for restaurants" (Step in the plan), home (phys. location), relaxing place (soc. location), evening (phys. time), relaxing time (soc. time), benevolence (PrioV, their two most important one), (0.6, 0.7) (WT1-level, WT2-level), size-based restriction (active Norms)]*. To ensure that guest have a table available at the restaurant, they select the achievement goal to have a table reserved which they pursue until they achieved it. A more interesting case in light of the size-based restriction, where we also see that actions of other groups affect agents is to consider a group of regular guests that come every Friday at the same time. They have their table reserved. This means that they are not affected by the new norm, and no external trigger happened. Now in the 'prepare the next step' plan, their 'next action' in the plan: 'order a beer' is not available anymore, as the restaurant owner cut the

beer in reaction to the new norm. This means that the regular guests are now in the 'find alternative action' step. While querying the action database, the guest finds the action: 'order a wine'. It has the same pre-conditions as the action to order beer: being at the table, having a waitress to take the order. Also, the post-conditions are similar. *Order_wine (table, have_waitress) = {{sit at table, waitress ready to take order}, {have wine, added wine to bill}} drink(wine) = {having pleasure}.* This action is not forbidden by the norm, and thus will be executed.

**We can now clearly see the short comings of other approaches**. For example: when the restaurant is fully booked and the regular guests bring one more friend this time. BOID [3] would handle such a conflict based on its agent types that make static decisions. [21, 22] reflect norm compliance decision with a utility function. The problem, however, remains: the decision is always the same, because the utility function will always give the same result (which might be that adhering to the norm has more utility than violating it). Reality is much more situated. Given the regularity of the guests, the restaurant owner might be more likely to let them in compared to novel guests.

Another issue that needs to be dealt with is the reaction of other agents to the agent's behavior. For example, the guests are going to react to the cheaper use of ingredients by the restaurant owner. Some are ok, others are not. While this is something to happen most likely in reality, it cannot be modeled by existing architectures. BOID [3], as well as [21, 22], have set responses, e.g., once the utility is defined for the use of cheaper ingredients, it will stay the same over the course of the simulation. Some guests might be fine if the norm is violated and a few more guests are in the restaurant than allowed. However, other guests might not be fine with that. Values play a crucial role here, and the current situation, e.g. how often did the restaurant owner violate the norm before.

## 5    Discussion and Related Work

We presented our architecture for norm deliberation that encompasses the motivational aspects of norms, their importance and full impact on planning and action selection. While existing norm models [1, 4, 11, 21, 22] have made great strides, they have several reasons why they are not suitable for our norm deliberation.

First, approaches like BOID [3] do not allow for the role of values nor context sensitivity. Neither does norm importance while planning and whether or not violating a norm play any role. Our architecture embraces the importance that values and current drive satisfaction play in making situated decisions and give different importance to different norms to allow for handling potential norm conflicts differently. This is contrary to BOID [3], where potential norm violations are always solved in fixed manner with pre-defined agent types always making static decisions, regardless of the situation.

Second, while utility-based approaches, such as EMIL-A [1, 4] and the work by [21, 22], are solving some of these issues they still do incorporate context sensi-

tivity for agents Whether to violate a norm or not is based on a utility function [22]. Then, the action (compliance or non-compliance) which provides more utility is chosen.

Nonetheless, utility functions have drawbacks, as they only work in static environments, and will always have the same outcome [11]. This static decision making is problematic, because we showed in our example discussion the dynamic nature of norm modeling, and the situatedness of decisions that influence normative behavior. Sometimes it might be more beneficial to adhere to a norm, while in other case it might be more beneficial to violate that norm. Defining a utility function for every possibility is not feasible. Such a function also needs to be modified and extended when an existing norm is modified or a new norm is added. Here, our value-based approach is more suitable, as our agents make decisions based on how the norm aligns with their values, and when an new norm is added, they just have to check the values connected to this new norm. Then they can react dynamically in every situation.

Third, while we use watertanks and values in our decision making based on [11, 16], we do not adopt the same approach as in [11] because we want explicit planning and reasoning with norms, as in their approach norms are only implicitly given by the actions. This makes it hard to modify a norm, because every action has to be inspected to see if it is potentially impacted by the change of the norm and the subsequent given needs satisfaction has to change. Another aspect where we deviate is by including explicit reasoning norm violations.

Because of the planning capabilities required by our agents, we see plan patterns as very useful here, because the agents only have to plan from landmark to landmark. Furthermore, we want to highlight that while we used a legal norm as an example norm, our architecture is also suitable for other types of norms, such as social norms or moral norms. For social norms for example, we can simply add a watertank for conformity. Conformity hear means to do what one's friends are doing [11]. In our immediate future work, we are going to implement our proposed agent architecture.

## References

1. Andrighetto, G., Campennì, M., Conte, R., Paolucci, M.: On the immergence of norms: a normative agent architecture. In: In Proceedings of AAAI Symposium, Social and Organizational Aspects of Intelligence, Washington DC. Citeseer (2007)
2. Brennan, G., Eriksson, L., Goodin, R.E., Southwood, N.: Explaining Norms. Oxford University Press, Oxford, Explaining Norms (2013)
3. Broersen, J., Dastani, M., Hulstijn, J., Huang, Z., van der Torre, L.: The boid architecture: conflicts between beliefs, obligations, intentions and desires. In: Proceedings of the Fifth International Conference on Autonomous Agents, pp. 9–16. Association for Computing Machinery, New York, NY, USA (2001)
4. Campenní, M., Andrighetto, G., Cecconi, F., Conte, R.: Normal= normative? The role of intelligent agents in norm innovation. Mind Soc. **8**(2), 153–172 (2009)

5. Castelfranchi, C., Dignum, F., Jonker, C.M., Treur, J.: Deliberative normative agents: Principles and architecture. In: Jennings, N.R., Lespérance, Y. (eds.) Intelligent Agents VI LNAI 1757, pp. 364–378. Springer (2000)
6. Crawford, S.E., Ostrom, E.: A grammar of institutions. Am. Polit. Sci. Rev. **89**(3), 582–600 (1995)
7. Dechesne, F., Di Tosto, G., Dignum, V., Dignum, F.: No smoking here: values, norms and culture in multi-agent systems. AI Law **21**(1), 79–107 (2013)
8. Di Tosto, G., Dignum, F.: Simulating social behaviour implementing agents endowed with values and drives. In: MABS, pp. 1–12. Springer (2012)
9. Dignum, F.: Autonomous agents with norms. AI Law **7**(1), 69–79 (1999)
10. Dignum, F.: Interactions as social practices: towards a formalization (Sept 2018). https://arxiv.org/abs/1809.08751v1
11. Dignum, F. (ed.): Social Simulation for a Crisis: Results and Lessons from Simulating the COVID-19 Crisis. Springer International Publishing, Cham (2021)
12. Dörner, D., Gerdes, J., Mayer, M., Misra, S.: A simulation of cognitive and emotional effects of overcrowding. In: Proceedings of the Seventh International Conference on Cognitive Modeling, pp. 92–98. Edizioni Goliardiche Triest, Italy (2006)
13. Duff, S., Harland, J., Thangarajah, J.: On proactivity and maintenance goals. In: Proceedings of the Fifth International Joint Conference on Autonomous Agents and Multiagent Systems, pp. 1033–1040 (2006)
14. Fikes, R.E., Nilsson, N.J.: Strips: a new approach to the application of theorem proving to problem solving. Artif. Intell. **2**(3–4), 189–208 (1971)
15. Gilbert, N., Ahrweiler, P., Barbrook-Johnson, P., Narasimhan, K.P., Wilkinson, H.: Computational modelling of public policy: reflections on practice. J. Artif. Soc. Soc. Simul. **21**(1), 14 (2018)
16. Heidari, S., Jensen, M., Dignum, F.: Simulations with values. In: Advances in Social Simulation: Looking in the Mirror, pp. 201–215. Springer, Cham (2020)
17. Kammler, C., Dignum, F., Wijermans, N., Lindgren, H.: Changing perspectives: adaptable interpretations of norms for agents. In: Van Dam, K.H., Verstaevel, N. (eds.) Multi-Agent-Based Simulation XXII, pp. 139–152. Springer, Cham (2022)
18. Kammler, C., Mellema, R., Dignum, F.: Agents dealing with norms and regulations. In: MABS, pp. 134–146. Springer (2023)
19. Mellema, R., Jensen, M., Dignum, F.: Social rules for agent systems. In: Aler Tubella, A., Cranefield, S., Frantz, C., Meneguzzi, F., Vasconcelos, W. (eds.) COINE XIII, pp. 175–180. Springer, Cham (2021)
20. Orkin, J.: Three states and a plan: the AI of fear. In: Game Developers Conference, vol. 2006, pp. 1–18. CMP Game Group SanJose, California (2006)
21. Panagiotidi, S., Alvarez-Napagao, S., Vázquez-Salceda, J.: Towards the norm-aware agent: bridging the gap between deontic specifications and practical mechanisms for norm monitoring and norm-aware planning. In: Balke, T., et al. (eds.) COIN, LNCS, vol. 8386, pp. 346–363. Springer, Cham (2014)
22. Panagiotidi, S., Vázquez-Salceda, J., Dignum, F.: Reasoning over norm compliance via planning. In: COIN, pp. 35–52. Springer (2012)
23. Schwartz, S.H.: An overview of the Schwartz theory of basic values. Online Read. Psychol. Culture **2**(1), 2307-0919 (2012)

# Validity Assessment of Uncertain Infection Indicators Using Virtual Artificial Society Model

**Yuki Misu and Shingo Takahashi**

**Abstract** Effective reproduction number is one of the indicators used to monitor the epidemic of an infectious disease. To calculate the effective reproduction number, it is necessary to know the time and route of infection of all infected people. However, since these cannot be observed in the real world, the number of new positives is used for estimation. In this paper, we focus on the uncertainty in estimating the effective reproduction number and propose a method to evaluate the impact of uncertainty in the real world using a virtual artificial society model. First, a model representing the infection situation is constructed, and the values of the effective reproduction number by definition are calculated in the model. It is possible to evaluate the validity of the estimated effective reproduction number by comparing the two calculated values of the effective reproduction number in the real world and the model. This method could replace unobservable uncertainty in estimation in the real world with "observable uncertainty in the model." Experiments are conducted to analyze the uncertainty of the rate of asymptomatically infected people and human flow. Finally, we discuss the results of experiments and their adaptability of this method to other fields.

**Keywords** Effective reproduction number · Validity assessment · Virtual artificial society

## 1 Introduction

More than 2 years have passed since the outbreak of COVID-19, but the epidemic has not yet ended in many countries around the world. The epidemic status can be assessed using various indicators. In most cases, data on actual infections are required to calculate these indicators. For example, the effective reproduction number ($R_t$) is

---

Y. Misu · S. Takahashi (✉)
Waseda University, Shinjuku-Ku, Tokyo 169-8555, Japan
e-mail: shingo@waseda.jp

Y. Misu
e-mail: arsenal0712@akane.waseda.jp

207

one of the typical indicators of epidemic status. The effective reproduction number means "the number of secondary infections caused by one infected person at a certain time $t$, under certain measures. However, the effective reproduction number cannot be accurately determined because infected people and their infection status cannot be identified sufficiently. Hence we need a new approach to estimating such indicators under uncertainty where enough data cannot be essentially obtained.

The effective reproduction number is usually defined by equation (1) [1, 2].

$$R_t = \int_0^\infty A(t, \tau) d\tau \tag{1}$$

The formula $A(t, \tau)$ expresses the secondary infection rate of infected individuals of infection age $\tau$ at time $t$. By including time $t$ as a factor, the indicator can represent changes in the infection rate due to virus mutation and infection control measures. In order to calculate the effective reproduction number according to equation (1), it is necessary to observe all infected people and their infection routes. It is not possible to observe all of them in the real world. Hence we need to estimate the effective reproduction number. For example, the following equation (2) is an estimation method using the number of new positive cases by the Japanese government [1, 3].

$$R_t = \left( \frac{J_{k+1}}{J_k} \right)^{\frac{\mu}{\Delta t}} \tag{2}$$

$J_{k+1}$ is the number of new positive cases in the last 7 days, $J_k$ is the number of new positive cases in the preceding 7 days, $\mu$ is the mean generation time ($= 5$ days), and $\Delta t$ is the reporting interval ($= 7$ days). Generation time refers to the period between the infection of the source and the infection of the secondary infected people.

Though the number of new positive cases should represent the number of all infected people, it can be estimated only as the number of patients diagnosed and reported. The average generation time is estimated to be 5 days, but this may change depending on the strain of the virus. Because of these uncertainties, it is essentially impossible to confirm if effective reproduction numbers are valid.

The purpose of this paper is to propose a methodology to evaluate the validity of indicators such as the effective reproduction number using a virtual artificial society model and agent-based simulation. Using a virtual artificial society, we can observe who are actually infected and calculate how many people are infected in the model.

## 2 Methodology Proposed for Evaluation of the Validity of Indicators

As mentioned above, we cannot observe all infected people and infection routes in the real world. On the other hand, in a virtual artificial society model using agent-based modeling, the infection situation can be "observed" in the model. Hence in the virtual artificial society, the effective reproduction number can be calculated based on its definition from the observation of the infection situation among agents in the model. This calculated value of the effective reproduction number must be considered as "true in the model" because the way of calculating the value satisfies the definition. At the same time, it is also possible to apply this estimation method mentioned above to "estimate" the effective reproduction number in the model as well. Thus effective reproduction numbers can be obtained in two ways: observation in the model and estimation by using the estimation method in the model. The results from the two ways can be different because of the different uncertainties involved in the real world and in the model. By analyzing these discrepancies, it is possible to evaluate the validity of indicators under uncertainties. This paper presents a method for evaluating the validity of uncertain indicators such as effective reproduction number using a virtual artificial society.

In the virtual artificial society model, infection situations can be actually observed in interactions between agents. Then secondary infection rates that are necessary for calculating effective reproduction numbers can be calculated in the model. Using the secondary infection rate calculated in the model, we can also calculate effective reproduction number, which we refer to as "measured $R_t$" in this paper.

By expressing the number of new positive cases used as an input value in the current estimation method in the model, the effective reproduction number using the estimation method can also be calculated. This calculated value is referred to as "estimated $R_t$" in this paper. The validity of the estimation method is verified by comparing the "measured $R_t$" and the "estimated $R_t$" in terms of the Euclidean distance (Fig. 1) [4].

The way of our proposed methodology of evaluating the validity of the effective reproduction numbers can be essentially applied to other indicators.

## 3 Overview of Virtual Artificial Society Model

### 3.1 Entities, State Variables and Scales

A virtual artificial society model in this paper should be built with sufficient components and interactions among agents to evaluate the validity of the effective reproduction numbers. Hence based on related models that were already validated [5], the virtual artificial model in this paper consists of behavior model, disease transition model, infection model, and diagnosis model. The model includes students, office

**Fig. 1** Overview of the validity assessment of the effective reproduction number using the virtual artificial society model

workers, housewives, medical workers, and the elderly as decision-making agents, and homes, schools, workplaces, hospitals, and commercial facilities as spots where agents visit, and are infected.

**Environment model**. The environment model represents two adjacent towns using a lattice graph. The model includes a home for each agent (*home*), a school (*school*) and a workplace (*work*) that are unique to each town, and a hospital (*hospital*) and a commercial facility (*amuse*) that are shared facilities between the two towns.

Each of these spots has a specific size, and *home* is divided into two types according to the agent's attributes. Other spots are places that agents visit and stay in during their daily activities, and the spots to visit are determined by the attributes of each agent and the town to which the agent belongs.

The size of *home* and *school* are set to the same values as in the previous study by Epstein et al. [5] in order to allow the same number of agents to visit. The values of *work*, *hospital*, and *amuse* and their sizes are calibrated to represent the stylized-facts as data on the accumulated number of positive cases in Shinagawa Ward from July 1 to August 31, 2020 [6]. This city was chosen as a typical case that has similar characteristics with the model to reproduce the approximate shape of the number of positive cases in the early stages of the spread of infection and the value of the effective reproduction number. Hence the proposed methodology could be applied to any city we would like to look into by building another model of that city (Table 1).

**Agent model**. Each agent has the parameters and elements shown in below, and follows the behavior model, disease transition model, infection model, and diagnosis model. The detailed interaction and emergence are described in the design concept.

**Table 1** Size and number of each spot

| Spot | | Size | Number of spots |
|---|---|---|---|
| *Home* | *Home_typeA* | $2 \times 2$ | 200 |
| | *Home_typeB* | $1 \times 2$ | 160 |
| *School(A/B)* | | $21 \times 11$ | 2 |
| *Work(A/B)* | | $10 \times 10$ | 2 |
| *Hospital* | | $30 \times 30$ | 1 |
| *Amuse* | | $14 \times 14$ | 1 |

Agent's parameters and elements

- *Id* 0 ~ 1119
- *group* {*st, sl, hw, hc, ag*}
- *city* A/B
- *pathology* {*susceptible, early, later, onset, heavy, recovered, death*}
- *infected* True/False
- *infectious* True/False
- *hospitalized* True/False
- *stay_at_home* True/False
- *infect_day* 0 ~ until recovery or death
- *source_Id* 0 ~ 1119
- *false_negative* True/False
- *infection_spot* {*home, school(A/B), work(A/B), hospitalamuse*}
- *behavior_after_onset*{*daily, hospitalized, stay_at_home*}
- *work_another_city* True/False
- *amuse_visit* True/False

*Agent type.* Each agent is defined by its attributes and the type of town to which it belongs. The types of attributes and the number of agents for each attribute are based on the previous study by Kurahashi et al. [7]. There are five types of attributes (student, salaried worker, housewife, health care worker, and aged) and two types of towns (A and B). Table 2 shows the number of people of each attribute in one town. The number of people with each attribute is the same in each town. The spots that the agents visit and stay are different according to their attributes and the town they belong.

**Table 2** Number of agents for each attribute

| Attribute | Number of agents |
|---|---|
| Student (*st*) | 200 |
| Salaried worker (*sl*) | 100 |
| Housewife (*hw*) | 95 |
| Health care worker (*hw*) | 5 |
| Aged (*ag*) | 160 |

*Behavior model.* In the behavior model, a day is considered as two steps. Each agent acts according to its own rules of action in the first step (daytime) and stays in its home in the second step (night). The rules of action are determined by the agent's attributes, town, and pathology.

*Disease transition model.* The disease transition model represents the pathology of an infected agent, from incubation period to recovery or death, the number of days of each pathology, and the probability of change. It is assumed that agents in the recovery state acquire immunity and do not become infected again.

*Infection model.* The infection model is based on the previous work of Epstein et al. [5], in which infection occurs at a spot. Uninfected agents are infected if there are infectious agents in the Neumann neighborhood at the visited spot, according to the infection probability determined for each spot.

*Diagnosis model.* The diagnosis model in this study is based on post-infection behavioral data of people infected with COVID-19 and information on the accuracy of PCR tests. A certain percentage of the infected agents are tested to obtain diagnosis results based on the accuracy of the test. In our model, only infected agents are tested, so the results are positive or false-negative.

## *3.2 Process Overview*

The flow of the simulation experiment to evaluate the validity of the effective reproduction number in this paper is shown below.

1. Agents are initially generated at the start of the simulation. The initial states are set to: $pathology$: "$susceptible$", $infect\_day$: 0
   $infected, infectious, hospitalized, stay\_at\_home$: $False$
2. Initially infected agents are generated from among the agents that are diagnosed after the onset, as they are thought to influence the spread of the infection.
3. Agents visit and stay at spots according to the Behavior model. At each spot, contact is recorded and infection occurs according to the infection model.
4. To calculate the "measured $R_t$", the infected agent acquires the Id of the infection source. It also outputs the number of new positives necessary to calculate the "estimated $R_t$".
5. At the end of each day, the "measured $R_t$" for that day is calculated. The "measured $R_t$" is calculated by summing the average number of secondary infections for each infection age on that day, based on the infection ages of the infection sources obtained in step 4.
6. The number of days is updated by one for each after two steps elapsed, and the age of infection of the infected agent is also updated.
7. Assuming two and a half months, the simulation is terminated if the number of days reaches 74. Otherwise, the simulation returns to step 3.

8. After the simulation is completed, the "estimated $R_t$" is calculated using the number of new positive cases obtained in step 4.
9. The Euclidean distance between the "measured $R_t$" calculated in step 5 and the "estimated $R_t$" calculated in step 8 is calculated to evaluate the validity.

## 4 Design Concepts

### 4.1 Basic Principles

The model in this paper is composed of a behavior model based on Epstein's virtual social model, which has been often used in previous studies on infectious diseases. Disease transition model, infection model, and diagnosis model are constructed based on published data and research findings on COVID-19. The model is enough to analyze validity of effective reproduction numbers.

Since the true value of the effective reproduction number cannot be known in the real world, it is impossible to examine how accurate the estimation is because of the influence of uncertainty in the real world. A virtual artificial society model contains uncertainty only when designing the model. Hence we can eliminate uncertainty from the comparison of "measured $R_t$" and "estimated $R_t$" obtained from the model. In other words, if a model of infectious disease satisfies a kind of validity criterion on the model, then the model could provide a way to evaluate uncertainty on the validity of the effective reproduction number.

### 4.2 Emergence, Interaction, Stochasticity

In the transmission of infectious diseases, people are infected through their actions at the micro-level, and the situation of society as a whole is observed with various indicators describing the characteristics at the macro-level. This model represents the general emergent nature of infectious disease transmission in such a way that agents infect other neighboring agents through their daily activities, and macro-level indicators such as the effective reproduction number and the number of new positive cases are recognized in the society as a whole.

Agent interactions in this model represent contact and infection, and do not include elements such as information exchange, learning, and adaptation. The following sections describe the detailed interaction, emergence, and probability of introduction of the model. We built each model presented below to generate the infected people necessary to calculate the measured values and the new positive people necessary to calculate the estimated values.

**Behavior model**. The following are the rules of action for agents. The action flow of the agents is shown in Fig. 2.

**Fig. 2** The action flow of agents

- Hospitalized: Stay in *hospital* until recovery or death
- Stay at home: Stay in *home* until recovery or death
- Daily activities

  (*st*) Attend *school* in the town where the agent belongs.
  (*sl*) $P_{other}$ percentage of agents commute to *work* in another town.

  Others commute to *work* in the town where the agent belongs.
- (*hc*) Commute to *hospital*.
- (*hw*, *ag*) A certain number of people are randomly selected to visit *amuse*.

Agents who recovered are assumed to return to their daily activities. $P_{other}$ is set to 0.1, the same value as in the previous study by Kurahashi et al. [7]. The number of *hw* and *ag* who visit *amuse* is initially set to 175.

**Disease transition model**. The disease transition model constructed based on previous studies of COVID-19 [8–10] and empirical data [11, 12] is shown in Fig. 3. The colors of the pathological states in Fig. 3 correspond to the colors of the agents in each pathological state on the model.

The period from infection to disease onset (incubation period) and the period from disease onset to behavioral change are determined by statistical distribution based on previous studies. The method of determination for each period, including the date of testing and the date of behavior change, is shown below.

- *incubation_Period*

  Lognormal distribution with mean 5.6 days and standard deviation 2.3 days

- *onset_to_test*

**Fig. 3** Disease transition model

Weibull distribution with shape parameter 1.741 and scale parameter 8.573 [13]

- $infectious\_day$

  $incubation\_Period/2$

- $behavior\_change\_day$

  $incubation\_period + onset\_to\_diagnosis$

The probabilities of severe illness and death differ depending on the agent's attributes. According to the Ministry of Health, Labour and Welfare [12], probabilities of severe illness and mortality for each attribute based on are shown below.

- Probability of severe illness

  $st, sl, hw, hc$: 0.003   $ag$: 0.085

- Probability of mortality

  $st, sl, hw, hc$: 0.0006    $ag$: 0.057

**Infection model**. The probability of infection for each spot was set based on the previous study by Klompas et al. [14] as follows.

- $home$: [0.1, 0.4]
- $school, work, hospital, amuse$: 0.05

**Diagnosis model**. The flow of diagnosis model is shown in Fig. 4.

In the previous study by Kurahashi et al. [7], it was assumed that 50% of the agents who developed the disease would self-treat and perform their daily activities, and since this study also assumes daily activities after the disease onset, we use that value at the reference time. In addition, from a publication by the Nara Medical Association [15], the sensitivity of the test is set at 70% because the average rate of false-negative results by PCR testing is approximately 30%.

**Fig. 4** The flow of diagnosis model

According to the "Data on Infection Status" published by the Ministry of Health, Labour and Welfare [11], it is known that among those who tested positive by PCR test before July 8, 2020, the proportion of hospitalized people is 73%, and the proportion of people receiving treatment at home or overnight stays is 27%. Based on this data, we set the proportions of hospitalizations and stay at home after the test.

## 4.3 Observation

The effective reproduction number calculated in the model and the number of new positives to calculate the effective reproduction number from the estimation equation are obtained as outputs.

# 5 Analysis Results

## 5.1 Analysis 1: Scenario Analysis for the Rate of Asymptomatically Infected People

Infections occur according to the infection rate in the infection model, and tests are performed according to the test rate and sensitivity in the diagnosis model. In this paper, scenario analysis was conducted according to the test rate to examine the effect of the proportion of asymptomatically infected people (The test rate is defined as the rate of the number of agents whose $behavior\_after\_onset$ are $hospitalized$ or $stay\_at\_home$). The average results of 100 trials for each scenario are shown in Table 3. The validity is evaluated with the Euclidean distance between the "measured $R_t$" and the "estimated $R_t$."

A t-test on the mean of the Euclidean distance showed that there was a 1% significant difference between the test rates of 0.1 and 0.3 with all other scenarios, and a 5% significant difference between the test rates of 0.5 and 0.9. In other words, the

**Table 3** Results of analysis 1

| Test rate | Number of positive people | Number of infected people | Capture rate | Euclidean distance (validity) | |
|---|---|---|---|---|---|
| 0.1 | 40.6 | 679.5 | 0.055 | 11.791 | (Low) |
| 0.3 | 108.1 | 635.1 | 0.159 | 9.833 | |
| 0.5 | 178.8 | 612.2 | 0.260 | 8.463 | |
| 0.7 | 245.3 | 598.5 | 0.377 | 8.226 | |
| 0.9 | 226.6 | 447.6 | 0.467 | 8.037 | (High) |

higher the test rate, the higher the capture rate of infected people and the higher the validity of the effective reproduction number.

Since the difference in test rates in this model corresponds to the proportion of asymptomatic infected people in the real world, it can be linked to the real-world phenomenon as follows: "In a situation where the number of asymptomatic infected people is high, the capture rate of infected people decreases and the relevance of the effective reproduction number is also low."

## 5.2   Analysis 2: Scenario Analysis for Human Flow

Next, a scenario analysis was conducted according to the number of people visiting a commercial facility per day in the model to examine the impact of human flow on a city. Table 4 shows the average of 100 trials for the number of positive people, the number of infected people, the capture rate, and Euclidean distance for each scenario.

No significant differences were found among the scenarios for each outcome measure. This shows that in our model the restriction policy of human flow would not directly lead to a decrease in the number of infections or an increase in the capture rate.

**Table 4** Results of analysis 2

| Number of visitors to *amuse* | Number of positive people | Number of infected people | Capture rate | Euclidean distance |
|---|---|---|---|---|
| 51 | 154.4 | 506.4 | 0.294 | 8.992 |
| 102 | 152.5 | 501.8 | 0.283 | 8.704 |
| 153 | 151.9 | 524.7 | 0.271 | 8.516 |

## 6   Discussion

In our method, the infection situation that is "unobservable" in the real world is made "observable" in the model, and the uncertainty of the situation in the estimation method in the real world is transferred to the uncertainty in the design of the model. This makes it possible to compare the "measured" and "estimated" effective reproduction numbers calculated in the model on the same scale.

In the analysis using the model, real-world uncertainties were treated as scenarios, and their effects on the validity of the effective reproduction number were evaluated to a certain extent. In this paper, Euclidean distance was used to compare the measured and estimated values. If differences in other types of data such as time-series ones should be evaluated, other indices would be required to quantitatively evaluate them.

The level of uncertainty given as a real-world element can be also represented in the model. For example, the impact of a delay in reporting data on the number of observed positives can be considered if it is represented in the model. Finally, we notice that the proposed method in this paper is general enough, and applicable to other cases even if the way of estimating the indicator in the real world is changed.

## 7   Conclusion

In this paper, we proposed a method that enables the comparison of "measured values" and "estimated values" of the effective reproduction number in a model using a virtual artificial society model, and evaluate the impact of uncertainty in the real-world situation on the validity of the effective reproduction number. The results of the scenario analysis indicate that the validity of the effective reproduction number decreases as the proportion of infected people who are not tested, such as asymptomatic infected people, increases. The results also suggested that increases or decreases in the flow of people in small-town units did not affect the infection situation.

The method of this paper can be applied to other situations of infectious disease epidemics by changing the model involved in the epidemic. The modeling concept of transferring the uncertainty of unobservable data to the uncertainty of the model setting also has potential for use in other cases such as the evaluation of economic indicators whose true values cannot be observed in the real world, and in the analysis of alternative data.

# References

1. Nishiura, H.: Effective reproduction number and surroundings. https://live2.nicovideo.jp/watch/lv325833316 (2021). Accessed 27 Jan 2021 (in Japanese)
2. Fraser, C: Estimating individual and household reproduction numbers in an emerging epidemic. PLoS ONE **2**(8), e758 (2007)
3. TOYOKEIZAI ONLINE: Status of COVID-19 infection in Japan. https://toyokeizai.net/sp/visual/tko/covid19/ (2021). Accessed 27 Jan 2021 (in Japanese)
4. Yoshikawa, T., et al.: A Study on Similarity Search for Very Long Time-Series Data. The Institute of Electronics, Information and Communication Engineers (2007) (in Japanese)
5. Epstein, J.M.: Studies in Agent-Based Computational Modeling, Generative Social Science. Princeton University Press (2007)
6. TOKYO Metropolitan Government COVID-19 Information Website: Positive cases (by municipality). https://stopcovid19.metro.tokyo.lg.jp/cards/number-of-confirmed-cases-by-municipalities/ (2020). Accessed 27 Jan 2021 (in Japanese)
7. Kurahashi, S.: Estimating effectiveness of preventing measures for 2019 novel coronavirus diseases (COVID-19). In: Proceeding of Conference on Social Systems. The Society of Instrument and Control Engineers (2020) (in Japanese)
8. Linton, et al.: Incubation period and other epidemiological characteristics of 2019 novel coronavirus infections with right truncation: a statistical analysis of publicly available case data. J. Clin. Med. **9**(2), 538 (2020)
9. The Novel Coronavirus Pneumonia Emergency Response Epidemiology Team: The epidemiological characteristics of an outbreak of 2019 novel coronavirus diseases (COVID-19) China. China CDC Weekly **41**(2), 145–151 (2020). https://doi.org/10.3760/cma.j.issn.0254-6450.2020.02.003(2020)
10. WHO-China: Report of the WHO-China joint mission on coronavirus disease 2019 (COVID19). https://www.who.int/docs/defaultsource/coronaviruse/who-china-joint-missionon-covid-19-final-report.pdf (2020). Accessed 27 Jan 2021
11. Ministry of Health, Labour and Welfare: Data on infection status. https://www.mhlw.go.jp/content/10900000/000649235.pdf (2020). Accessed 27 Jan 2021 (in Japanese)
12. Ministry of Health, Labour and Welfare: 10 things to know about COVID-19 infections now. https://www.mhlw.go.jp/content/000689773.pdf (2020). Accessed 27 Jan 2021 (in Japanese)
13. Nishiura, H., et al.: Serial interval of novel coronavirus (COVID-19) infections. Int. J. Infect. **93**, 284–286 (2020)
14. Klompas, M., et al.: Airborne transmission of SARS-CoV-2: theoretical considerations and available evidence. J. Am. Med. Assoc. **324**(5), 441–442 (2020)
15. Nara Medical Association: Testing for COVID-19 infections. http://nara.med.or.jp/for_residents/11121/ (2021). Accessed 27 Jan 2021 (in Japanese)

# Social Identity and Social Influence

# First Step Towards a New Understanding of Radicalisation: Modeling Identity Fusion

**Mijke van den Hurk, Mark Dechesne, and Frank Dignum**

**Abstract** We want to understand in which circumstances identity fusion occurs. We propose a model in which individual needs and interactions between agents and their social environment come together. We argue the personal identity of an agent will fuse with a group when it has a high need for significance and he is member of a group providing a means to gain significance. Agents cannot join all groups to meet their needs, as agents need to have a social connection with the group and need to be accepted within the group. The model allows for multiple scenarios to occur. Agents with a need for significance not necessarily become fused and will find alternative ways to satisfy their need.

**Keywords** Identity fusion · Personal identity · Social identity · Agent-based modeling · Pro-group behaviour

## 1 Introduction

We want to understand in which circumstances identity fusion occurs. Identity fusion is the process in which an individual fuses with a group identity and "*the group comes to be regarded as functionally equivalent with the personal self*" [1]. Fused people retain the feeling of making autonomous choices about their behaviour, while their

M. van den Hurk (✉) · F. Dignum
Utrecht University, Utrecht, The Netherlands
e-mail: m.vandenhurk@uu.nl

M. van den Hurk
Dutch National Police, The Hague, The Netherlands

M. Dechesne
Universiteit Leiden, Leiden, The Netherlands

F. Dignum
Umeå University, Umeå, Sweden

CVUT Prague, Prague, Czech Republic

social identity with relation to the group stays salient too. The activation of both the personal and social identity leads to high motivation for performing actions benefiting the group while making personal sacrifices.

Our interest in this phenomenon is the link with violent extremism. The degree of fusion of an individual is a strong predictor of pro-group behaviour where people are willing to sacrifice personal costs in order to benefit one's group [2]. Also, a potential threat for the group identity means a threat towards the personal identity, and one wants to protect the group and therefore himself. Costs can be in terms of money or time, and, in extreme cases, the willingness to use violence or give one's life to protect one's group. Given the willing for self-sacrificing actions fused people tend to have, the question on how some people radicalise towards violent extremism cannot be answered without insights in the mechanism behind identity fusion.

Research on identities is complex since multiple variables on individual and social level come together. Individuals can fulfill their social needs by connecting with groups. Groups exists out of individuals itself, but also relate to other groups, based on what they found important. The interaction within and between the two levels make the process of identity fusion a complex mechanism. An agent-based model is a useful research method to analyse such dynamic and complex processes. It gives us the opportunity to define the interactions between individuals and groups, and study the development of a fused identity over time.

In this paper, we will describe an agent-based model representing a society with agents having social needs, namely the need for significance, i.e. the need to be acknowledged for one's actions, and the need for belonging, i.e. the need to connect with others. Agents become member of groups for which they develop a social identity. The groups represent real-life groups, such as colleagues, friends and political parties. Actions in the model are performed by agents to fulfill the personal needs, but also contribute to the vitality of their groups. Agents cannot join every group in the model, as they need to have a social connection with a group before becoming a member, representing the likelihood of people joining a group when they have a friend who is already a member. Also, agents need to be accepted within a group. We model the social identity as the overlap between the personal self and the group identity, and the personal identity as the collection all social identities [3]. We define a fused identity as an agent having one single dominant social identity and willing to make personal sacrifices in favor of the group, while maintaining need satisfaction. We can use the model to study in which scenarios fusion occurs, i.e. which combination of individual needs and available groups leads to a fused identity.

The rest of the paper is as follows. Section 2 contains related work. Section 3 will go into the process of identity fusion and the base of the model around identities is explained. Section 4 describes the model in more detail. We will conclude with a discussion about the model and future work.

## 2 Related Work

Our model is not the first ABM with social or personal identities. The model in [4] was based on the optimal distinctiveness theory, where agents' group choice depends on social identities of others. It lacked, however, a concept of a personal identity. Without a personal identity, we are not able to model identity fusion, i.e. the process of a personal and group identity becoming aligned. Prada et al. [5] showed how to integrate different personalities within a framework of group dynamics. Here, the relation with social identities is missing in order to model identity fusion. Finally, [6] shows with a ABM how crowd behaviour emerges based on the Social Comparison Theory. The theory could explain identity fusion with a group on short term, but as soon as the group disappears, people become defused. So, the model is not useful to explain the process of becoming and remaining fused for a longer time span.

## 3 Process of Identity Fusion

Identity fusion theory explains the concept of identity fusion, but is not explicit about how it emerges. The theory builds around the idea that the personal self and group identity can merge with each other, and both the self and the social identity become salient at the same time. Traditional research on social identities assumes personal and social identities are two separate concepts, i.e. the former describes a person as an individual and the latter an individual in social context [7]. The more the personal identity becomes salient, the less the social identity is active, and vice versa. However, with this theory we cannot explain how a personal identity merges with a group identity. We therefore propose to use the interaction between the concepts to explain the emergence of identity fusion. We will use the rest of this section to explain the core concepts of the model and how identity fusion emerges.

The model represents a society in which agents strive for need satisfaction. Needs are satisfied by performing actions and belonging to a group [8]. The groups represent real-life groups, such as family and friends, colleagues and teammates, and fellow members of a political party or religion. Groups We will describe the behaviour of the model by looking at an agent. Let's call him Bob. Bob works at an organisation helping refugees. He has a family and some friends, and likes to play soccer every Saturday.

*Social identity* When Bob started to play at the soccer team he belonged to that group and a social identity was formed. We define the corresponding *social identity* as that part of the group identity one identifies with. It is a collection of actions, i.e. "how should I behave" and attitudes, i.e. "what do I think of this". The social identity gives guidance in how to behave within the group. Bob knows skipping soccer practice is rejected by the team and that the team likes to have a drink after a match on Saturday.

*Personal identity* Bob belongs to multiple groups, i.e. his work, soccer team, family, etc., and therefore has multiple social identities. We define *the personal identity* as

the collection the social identities related to groups he is a member of. Each identity has a weight, representing the commitment towards that identity. The needs are used to determine which social identity should be prioritised over others. If Bob wants to belong to a group he joins the Saturday drinks with his teammates, and if he wants to be acknowledged by his team, he joins every practice to become better. If, however, his family is more important to him, he will chose family gatherings on Saturday over drinks with the team for more feeling of belonging.

*Fused identity* We can now define a *fused identity* as an agent having one main social identity which is salient at all times, being prioritised by the personal identity. Other social identities can exist, since we all have multiple identities such as gender and nationality, but they do not counteract on the main social identity. Furthermore, the fused person is willingly performing pro-group actions and is still able to fulfill personal needs. He is not forced to do certain actions, but chooses to behave as the group. We can analyse the degree of identity fusion of agents by measuring the number of social identities adopted in their personal self, the extent to which one social identity is prioritised over others and their level of need satisfaction.

*Needs* The needs of the agents are the need for significance, belonging and survival. We define *the need for significance* as wanting to perform actions that are approved of within a group and benefit the group but come with personal costs [9]. Bob can gain significance by participating in the match every Saturday or by performing well on his job, etc. The higher the personal sacrifices, the higher the increase in significance. Doing a disapproved action will make significance decrease. If Bob chooses to go to the bar with friends instead of joining soccer practice, his significance will gain because of his friends, but will also decrease as his soccer team condemns his choice. The level of significance is only affected by groups Bob is member of. Costs of actions will decrease the level of survival satisfaction. We define *the need for survival* as the level of personal resources. Bob's level of survival increases if he chooses not to perform any action, as in that case he will have time to eat, sleep, etc. The costs of an action decrease as the action is performed more and hence less significance can be gained. Bob is becoming better at soccer after some months, but people now know he will score at least once every match. He has to put more effort in and at least score two times in order to be acknowledged by his time as much as when he started to become better.

We define *the need for belonging* as the need to perform in groups, while being accepted within the group. The need can be satisfied by performing actions of a group of which one is a member. A high need for belonging will drive an agent to join a group with enough members. The bigger the group, the more the feeling of belonging can be met. The level of costs of the action is irrelevant. If Bob is tired and not playing very well during the match, he still participated within the team and his feeling of belonging is as high as when he performed better.

*Decline in significance* A decline in significance satisfaction makes the agents look for a means to gain it back and fusion might start. For instance, Bob is fired. He is feeling useless and underappreciated, i.e. his significance is lowered, and hence

his need for significance becomes high. With his free time, he decides to help his mother with daily courses and extra soccer training to become a better play helps him increase his need for significance.

Now imagine Bob likes his soccer team, but his team does not have the ambition to win, while Bob needs more recognition for his actions. Bob talks to a friend about this and his friend tells him they are looking for volunteers at a new political party. Bob agrees on the ideology of the party and starts working 3 three days for them. Bob is being appreciated for his time and he feels better. Meanwhile, he finds another job and has to reduce his work for the political party to one day a week. Bob is happy. In another reality Bob is not able to find a job. His feeling of significance drops after every rejection. Fortunately, his political party is happy with him. After a few months of volunteer work he is asked to join the campaign committee as his communication skills are great. He spends more and more time at the party's office and starts to make friends there. He feels appreciated and the party starts to predominate his life. Fusion starts to emerge.

What if Bob does not agree with the statements of the political party? The party being against immigration. Joining the party would not only counteract on his own beliefs, but also on his friends and family with whom he shares the same beliefs. If he remains pro-immigration, he will not be accepted in the party. Bob has weigh his options: does he want to change his belief about immigrants in order to gain significance, or find another means? We model these deliberation about identities according to the balance theory [10]. Bob decides his family is more important to him and does not become a member of the party or, alternatively, he starts to fight anti-immigration groups.

We described scenarios in which Bob was able to fulfill his need for significance. The availability of a group in which significance can be obtained was necessary for need satisfaction, but not necessarily leads to fusion. Fusion only emerged because he had a connection with a group, i.e. the political party, and this group was the only means to gain significance. He remained, however, committed to former social identities, as they did not counteract on his political identity.

While deciding about joining the anti-immigration party or not, he also takes his need for belonging into account. He loves his family, but they are busy with their lives and he does not see often. The political party has a lot of members doing activities regularly. This is an extra reason for him to join the party after all. It puts Bob in a difficult position during family gatherings. Although the party gives him significance, being with his family actually decreases it. Bob chooses to skip the family gathering to avoid the direct rejection of them. Bob now has less groups he can attend, and his social identity of the party becomes more predominant. As he becomes more prominent within the party and is less connected with other groups, his only way to gain significance is by putting in more personal effort for the party. His family will judge him even more, and he becomes detached from them. It accelerates the process of identity fusion even more as his personal identity has less social identities to chose actions from.

   We showed different scenarios in which a low level of significance could be increased, depending different circumstances. In only some of the scenarios this led to identity fusion. We want to verify our hypothesis by building an agent-based model and simulate the above scenarios. We expect fused agents to become isolated in the core of groups, where they are only surrounded by members of the same group, while remaining the option to connect with other groups.

## 4 An Agent-Based Model

We will use this section to explain our model in more detail. We will use the model from [11] as a base model.

### 4.1 Agents

Agents in the model have three need satisfaction levels $S_{sig}$ (significance), $S_{bel}$ (belonging) and $S_{sur}$ (survival) $\in (0, 1)$. Significance and belonging decline over time and actions need to be performed to let the levels increase, while survival increases over time, i.e. doing nothing gives the agent time to recharge. If a need is below threshold $\lambda$, action is required. The threshold values differ among agents, i.e. some will always have more urge for significance or belonging than others. Agents can choose between multiple actions from action collection $A$. An agent can only choose actions from groups he is a member of. Every action takes personal resources, representing time, money and effort, and makes the survival level drop. Significance is gained by choosing an action that has a positive effect on the group vitality and is not performed by others. The more personal resources are used, the more significance is obtained. Belonging is earned by performing actions in favour of the group with other members of the group. The higher the average group commitment of its members towards the group, the more belonging can be gained.

**Social identity** When an agent chooses an action of a group a social identity is formed. The social identity contains all actions associated with the group and corresponding attitudes. The social identity $S$ for group $i$ has a weight $w_{S_i}$, representing the commitment.

**Personal identity** The personal identity of the agent is the collection of social identities (Fig. 1a) and represented in an associative network (AN) (Fig. 1b). It consists out of relations between the agent, *Self*, his social identities, actions and commitment towards a social identity. The relation between *Self* and an action is drawn when the action is performed at least once. For now, we do not draw lines between social identities. Commitment for a new group starts at a random level. By default, the commitment towards social identities decreases over time. The weight is updated after the agent performed an action from the identity with the change in satisfaction.

(a) Personal identity, which is the cumulative overlap of the group identities one belongs to.

(b) Example of an AN. The *Self* is member of group 2 and 3. Action $a2$ is performed, hence the positive attitude between *Self* and $a2$.

**Fig. 1** Conceptual visualisations of the personal identity (**a**) and the relation between the self and other social identities in an associative network (**b**)

Thus, when an group provides actions resulting in an increase in need satisfaction the commitment towards that group strengthens. The higher the commitment towards a social identity, the more defining the associated group becomes for the agent and hence the more central in identity fusion.

*Maintaining a balanced personal identity* Adding a new social identity can potentially cause an imbalance, as the same action could be evaluated differently by different groups. Balance is determined by relating *Self*, the new action and the social identity to each other with their relative attitudes. We can speak of balance if there are no or two negative attitudes between three concepts. The personality identity determines which social identity potentially contributes to the highest need satisfaction.

## 4.2 Actions

Actions associated with the group identity tell agents what actions contribute to group vitality and therefore give significance and/or belonging. Each action has two variables: i.e. valence and cost. The costs are independent of the group in which the action is performed. Valence represents the attitude or affecting quality of an action within the group. The valence can be positive or negative and indicates how much significance an agent can gain by performing that action. Different groups can have different valence values for the same action.

*Intensity* Each action can be performed at different intensity levels. A higher level intensifies the costs, and, thus, the potential gain or loss in significance. The action *a contribute to political party* for example can be performed with low intensity or with high intensity. With low intensity we model an agent choosing to come to the annual general meeting. This will only cost him a little bit of his time, and showing up will give him a small gain in significance. However, choosing to volunteer once a

week by joining the board of the party will cost an agent more of his time, but as he is appreciated more for his actions, his gain in significance will also be higher. We model this by letting the agent perform action $a$ with a high intensity.

After performing an action at a certain intensity a few times, the costs decrease as the agent gets better at it. This also means a smaller increase in significance. The agent therefore has to increase his intensity over time to ensure a gain in significance. Agents can always choose actions at a level of intensity they performed before or with a +1 increase. A new action is always performed at intensity level 1.

## 4.3  Groups

We assume groups $G_1$ to $G_n$. A group consists of agents and has an identity, which is a collection of actions $A_{G_i}$ and attitudes towards those actions. We will relate these concepts in a associative network. The group identity is independent from the identities of its members. An aggregation function of personal identities can make group identities dynamic, but for simplicity we chose not to do so. Each group has vitality $v_{G_i}$ representing how well the group is doing on reaching group goals. The vitality decreases over time and, as with the needs and group commitments of agents, actions have to be performed to keep the vitality as high as possible. Each action in the group identity influences the group vitality.

We want a variety of group identities as we want to analyse which type of group identities leads to identity fusion. Variety is created by letting groups have actions which are not included in other identities or have opposed attitudes towards the same actions, i.e. counteracting actions, such as being pro- and anti-vaccination. Furthermore, we define three group characteristics. First of all, not all group identities contain actions which can be performed with others, but only individually. These actions will not contribute to belonging. Secondly, some groups set a limit to how much personal sacrifices can be put in performing an action. Actions of a group of friends playing a game or having a drink together cost little time or money and only small amounts of significance can be gained. A political party needs volunteers or board members, and a sports teams requires practice and, therefore, result in more significance. These groups are not beneficial for gaining a lot of significance. Finally, we define loose versus tied groups, which refers to the spreading $I$ of allowed intensities action can be performed with. The intervals are set around the most chosen intensity level of actions by its members, which means the interval moves over time. Loose groups have a wide spread, so agents are allowed to deviate from average group behaviour. Tied groups on the other hand force agents to choose intensities similar to those of other group members.

**Joining a new group** Agent can only become a member of a group if they meet two requirements: proximity and acceptance. The agent has to be in close proximity of a member of the potential new group, modelling the importance of having social connections to new groups before becoming a member. An agent can only join group

$G_i$ if another member of one of his already joined groups is also member of group $G_i$. Secondly, the agents needs to be accepted within the group. An agent is accepted when he has fits the intensity requirements of the group, i.e. he performed an action at intensity level $l$ and $l \in I_{G_i}$.

## 4.4  Model Behaviour

Agents are member of multiple groups from the start of the simulation. Depending on the type of group they are a member of, they are able to fulfill their needs of significance and belonging. At the start of the deliberation they determine if action is needed. If so, they select all the possible actions from their personal identity, i.e. the collection of actions from their social identities. They compute which action is most beneficial, where the satisfaction levels determine the urgency for one or another action. When the need for belonging is lower than the need for significance, actions that can be done together are preferred over actions accounting for significance. Each agent has a so-called *action schedule*. It is a sequence of the groups he is member of, and requires agents to perform an action with group $G_t$ at tick $t$. It represents daily life where one has to work from Monday to Friday, and play soccer at Saturday. One can skip the requirement, but this will lead to a penalty, i.e. a decrease in belonging.

Some agents will not be able to fulfill their needs within their existing groups. This depends on the combination of the groups they belong to, i.e. can significance and belonging be gained, of their own needs, i.e. a low versus high threshold for the needs, and the possibility of performing an action at a feasible level, i.e. does the group allow the agent to perform at an intensity level he can reach. The agent has the option to change groups if his needs stay below satisfaction levels for a number of ticks. He can choose to join a new group, with new actions, but only when the proximity and acceptance requirements are met. He needs to share a group with a member of the new group, and he must be allowed to perform the new actions at level 1. He can also join a group with actions already available in his personal identity, but the same requirements hold. When the group contains new actions he must ensure a balanced identity.

## 4.5  Emergence of Identity Fusion

With the above proposed model we can study in which circumstance an agent becomes fused. We defined identity fusion as the process in which an agent develops one dominant social identity. Other social identities can exists but they contain no actions counteracting on the main social identity. Furthermore, the agent still has the option to choose actions from other identities but prefers his main social identity.

An agent develops such a main social identity when the actions from the social identity are chosen over other identities, which is the result of the positive feedback

**Fig. 2** Different scenarios for an agent with a high need for significance. Scenario (**a**) shows a balance between identities. Fusion will only occur in scenario (**b**) and (**c**), while in (**d**) the agent has to look for another group

loop between need satisfaction and commitment towards a group and the motive for a balanced personal identity. Different scenarios can occur when an agent has a high need for significance, of which only some result in identity fusion, see Fig. 2.

First of all, the process starts when an agent has a high need for significance. It will motivate the agent to join a group where significance can be gained. The more resources an agents put in, the more significance is earned. This will result in a positive evaluation of the group and, thus, a high commitment towards the associated social identity. If the agent has multiple groups to gain significance, no emergence will occur (a). In (b) the agent can join group 1 where more significance can be earned than in group 2. He maintains cognitive balance. As he has only one group to gain significance, he has to choose actions with a high intensity level. It results in a higher commitment towards group 1 and fusion emerges, while remaining a commitment towards group 2. In scenario (c) the agent is member of group 2 associated with a counteracting action regarding group 1. He has to distance himself from members of group 1 to remain cognitive balanced. He has to choose actions with a high intensity level to gain significance to and overcome the rejection from group 1. It leaves the agent with less groups to choose from to satisfy his need and fusion emerges. Finally, an agent will not always be able to join a group where high levels of significance can be earned (d). If the agent is not surrounded by a member of such a group, he is forced to gain little amounts of significance from groups he can join. Also, when such a group exists in one's environment, the agent should be accepted within the group. If not, fusion will not emerge.

## 5  Discussion and Future Work

We build a model to study the process of identity fusion. We defined identity fusion as developing a dominant social identity which the personal identity is committed to. We expect identity fusion to start with an agent having a high need for significance, i.e. being acknowledge by one's group for sacrificing personal costs in favor of the group. An agent has the option to join a group providing a means to gain significance as long as he is in close proximity with a group and is accepted within the group. The need drives the agent to choose actions with a higher intensity level over time, sacrificing personal costs. As long as the agent has multiple groups to gain significance from

and can commit to a variety of social identities, fusion will not happen. An agent will develop a dominant social identity if it has only one group to gain significance from. The agent will stay connected with other groups as long as the social identities ensure cognitive balance. If, however, his main social identity contains counteracting actions the drive for cognitive balance enforces him to avoid other groups. Identity fusion is accelerated when less social identities are contained in the personal identity. The need for belonging on itself will not cause identity fusion, but will boost the process of fusion in combination with a high need for significance.

With this model, we show how to integrate personal and social identities for which we have taken multiple theories from sociology and psychology into account. Furthermore, we show how both needs on an individual level and groups on a social level can be put together. We also showed that the combination of the theories can be build into one mechanism in which agents develop social identities in general, and identity fusion might emerge.

The model can be extended on multiple levels. We modeled identities as the collection of actions and attitudes towards those actions. Values should be included too as they are an important concept regarding group identities and ideologies. Research on identity fusion shows that the strongest fusion emerges within groups containing values in their identity [1]. Furthermore, we only took social identities into account of surrounding groups. However, we all have an identity based on gender and nationality. These are not primarily based on social groups but do play an important role in our personal identity. Therefore, these should be integrate to. Finally, we left out the relation between groups. Some groups have a negative attitude towards other groups, such as opponent sport teams or political parties. Integrating this relation gives agents the possibility to reason about groups and members of the groups. We can use this to model ingroup superiority and outgroup inferiority.

Although the model is a simplification, it creates a complex model with multiple variables interacting with each other. The next step would be implementing the model and simulate the mechanism of developing social identities. We can validate the model by proving all variables are needed for identity fusion to happen, and analyse the scenarios in which identity fusion occurs.

# References

1. Swann, Jr, W.B., et al.: Identity fusion and self-sacrifice: arousal as a catalyst of pro-group fighting, dying, and helping behavior. J. Person. Soc. Psychol. **99**(5), 824 (2010)
2. Swann, Jr, W.B., Buhrmester, M.D.: Identity fusion. Curr. Directions Psychol. Sci. **24**(1), 52–57 (2015)
3. Sharma, S., Sharma, M.: Self, social identity and psychological well-being. Psychol. Stud. **55**(2), 118–136 (2010)

4. Smaldino, P., et al.: An agent-based model of social identity dynamics. J. Artif. Soc. Soc. Simul. **15**(4), 7 (2012)
5. Prada, R., Camilo, J., Nunes, M.A.: Introducing personality into team dynamics, 667–672. ECAI 2010. IOS Press
6. Fridman, N., Kaminka, G.A.: Modeling pedestrian crowd behavior based on a cognitive model of social comparison theory. Comput. Math. Org. Theor. **16**(4), 348–372 (2010)
7. Tajfel, H, et al.: An integrative theory of intergroup conflict. Org. Identity Reader **56**(65), 9780203505984-16 (1979)
8. Dörner, D., et al.: A simulation of cognitive and emotional effects of overcrowding. In: Proceedings of the Seventh International Conference on Cognitive Modeling. Edizioni Goliardiche, Triest, Italy (2006)
9. Kruglanski, A.W., et al.: The psychology of radicalization and deradicalization: how significance quest impacts violent extremism. Polit. Psychol. **35**, 69–93 (2014)
10. Greenwald, A.G., et al.: A unified theory of implicit attitudes, stereotypes, self-esteem, and self-concept. Psychol. Rev. **109**(1), 3 (2002)
11. van den Hurk, M., Dignum, F.: Towards fundamental models of radicalization. In: Conference of the European Social Simulation Association. Springer, Cham (2019)

# HUM-e: An Emotive-Socio-cognitive Agent Architecture for Representing Human Decision-Making in Anxiogenic Contexts

**Patrycja Antosz** ⓘ**, F. LeRon Shults** ⓘ**, Ivan Puga-Gonzalez** ⓘ**, and Timo Szczepanska** ⓘ

**Abstract** This paper introduces HUM-e emotive-socio-cognitive decision-making architecture of human agents. The HUM-e is an extension of the HUMAT architecture representing socially influenced decision-making. The new architecture was designed to be used in contexts where fear and social identity play significant roles in attitude formation. Crucially, we propose that fear changes the goal of information exchange between interlocutors and influences the persuasiveness of information sources in face-to-face communication. The development of HUM-e was theoretically inspired by emotional contagion theory, identity fusion theory, and social identity theory.

**Keywords** Emotional contagion · Social identity · HUMAT · HUM-e · Decision-making · Attitude formation

## 1 Introduction

The goal of this paper is to introduce the HUM-e emotive-socio-cognitive decision-making architecture of human agents. The next section briefly grounds the architecture in previous scholar work. Subsequently, we describe HUM-e with the use of the ODD protocol [1].

P. Antosz (✉) · F. L. Shults · I. Puga-Gonzalez
NORCE Norwegian Research Centre AS, 5008 Bergen, Norway
e-mail: paan@norceresearch.com

T. Szczepanska
UiT The Arctic University of Norway, 9019 Tromsö, Norway

## 2 The Need for a New Decision-Making Architecture

Early approaches to modeling human behaviors focused mostly or solely on human rationality, simply assuming the validity of rational choice theory. Over the last few decades, perfect rationality has been robustly challenged by empirical findings and theoretical contributions [2]. Consequently, numerous architectures representing human decision-making processes have been developed (for a review, see [3]). Users of agent-based representations of human decision-making devoid of the emotional aspects found them lacking and began working to change the status quo. One group of scholars have developed new models explicitly focused on accounting for the role of affect in human attitudes and behaviors. Their models or architectures are based on theories such as the appraisal theory of emotions [4, 5] and Bayesian affect control theory [6]. Inspiring as computational theories, the more general representations pose significant challenges when calibrating to real-life cases. Another group of scholars have altered established cognitive ABM architectures to incorporate the affective dimension, e.g., ACT [7] or CLARION [8], EBDI [9, 10]. This work belongs to the second group.

Here, we described our work extending the HUMAT socio-cognitive agent architecture with fear. The driving force behind our work was the need to represent human attitude formation and decision-making in a novel, COVID-19 context. Quantitative data abundantly clear about the importance of fear for COVID-19 compliance [11–13] meant that computational representations of the phenomenon should incorporate the fear factor. Acknowledging the role fear plays in attitude formation in addition to the socio-cognitive motives present in HUMAT (i.e., experiential needs, social needs, and values), the HUM-e architecture explicitly specifies how fear is synergised with the cognitive and social dimensions of information processing and decision-making. The specification of assumptions related to fear and social identity was rooted in emotional contagion theory, identity fusion theory, and social identity theory (see Sect. 3.2.1).

## 3 HUM-e

### 3.1 Overview

#### 3.1.1 Purpose and Patterns

The purpose of the HUM-e emotive-socio-cognitive architecture is to represent the processes of attitude formation and decision-making when at least one choice alternative is perceived as harmful to one's health or well-being. The HUM-es exchange information to form attitudes toward choice alternatives ($a_1, \ldots, a_n \in A$), and subsequently choose the preferred alternative. The architecture can be extended to allow for a relevant number of alternative choices.

### 3.1.2  Entities, State Variables, and Scales

Agents represent individuals in the modelled population. Each agent has a social identity designating personal group affiliation. Agents identify themselves as members of one of three groups ($g_1, g_2, g_3 \in G$). Group membership describes an aspect of identity that is particularly important in the context of the choice the agents make (e.g., authoritarian, centre, or libertarian political views). Communication between in-group members (e.g., two agents who identify as $g_1$) always aims at reaching consensus, independent of speakers' agreement about the chosen alternative. Communication between members of different non-antagonistic groups (e.g., $g_1$ and $g_2$ or $g_2$ and $g_3$) builds consensus only of the speakers agree. Otherwise, when the speakers choose different alternatives, communication results in ignoring the interlocutor. Last, communication between antagonistic group members $g_1$ and $g_3$ results either in repulsion (if speakers advocate different alternatives) or in consensus (if speakers agree about their choices; see *Interaction* for details).

Each agent is characterized by a set of motives that are relevant in the context of making a particular decision. Motives can be related to experiences, values, health-concerns, and social inclusion ($e_1, v_1, h_1, s_1 \in M$). The architecture can be extended to allow for any number of choice-relevant motives. The agents vary with respect to how choice alternatives ($a_1$ and $a_2$) evaluate the motives. Evaluation of a choice alternative ($E_{a,m,j}^{t_n} \in \langle -1, 1 \rangle$ is a sum of evaluations of individual motives:

$$E_{a,j}^{t_n} = \frac{\sum_{m=1}^{4} E_{a,m,j}^{t_n}}{4} = \frac{E_{a,e_1,j}^{t_n} + E_{a,v_1,j}^{t_n} + E_{a,h_1,j}^{t_n} + E_{a,s_1,j}^{t_n}}{4} \tag{1}$$

$E_{a,m,j}^{t_n}$—evaluation of choice alternative $a$ with respect to motive $m$ for HUM-e $j$ at time $t_n$.

Choice alternative can be evaluated:

- negatively: $-1 \leq E_{i,m,j}^{t_n} < 0$,
- neutrally: $E_{i,m,j}^{t_n} = 0$
- positively: $0 < E_{i,m,j}^{t_n} \leq 1$.

In an implementation of HUM-e, calibration of initial evaluations of motives for choice alternatives can be linked to in-group beliefs and values.

Links between the agents enable communication acts i.e., sharing information about the expected consequences of each choice alternative of ($a_1$ and $a_2$) for the relevant motives. Information exchange between agents can take two forms: signalling and inquiring (see *Information exchange* for details). One time step of the model is an abstract unit and represents the trigger for discrete events.

### 3.1.3   Process Overview and Scheduling

As depicted in Fig. 1, HUM-e consists of two main processes executed at each time step: (1) information exchange in social networks, and (2) attitude formation (see Eq. 1).

**Information exchange**

As the agents are making up their mind about what alternative to choose, they reflect on pros and cons of both choice alternatives and create an expectation on how satisfying each decision will be for them. For each choice, if at least one motive is evaluated positively and at least one negatively, the agent experiences an unpleasant state of cognitive dissonance. This ambiguity (see *Initial choice* section of *Initialization* for details) occurs because the evaluated alternative is neither uniformly good nor uniformly bad for all relevant motives, but rather has both pros and cons. Consequently, the agent faces a dilemma and, if the dissonance is strong enough, engages in information exchange to reduce it and maintain cognitive consistency. The agent strives to change either its own beliefs or the beliefs of others in the agent's social networks, so that the preferred alternative becomes more internally consistent. The agent achieves this either by signalling to others (attempting to convince others that they should do what the persuading agent chose to do) or by inquiring of others (asking for advice).

The dissonance reduction strategy implemented by the agent depends on the configuration of pros and cons that the preferred alternative evokes. The agent can face two types of dilemmas: social and non-social, which lead to two different types of information exchange: signalling and inquiring, respectively. A social dilemma occurs when the preferred alternative yields negative evaluation of the social inclusion motive and a positive evaluation of any of the other motives. It corresponds to a situation in which the agent is convinced that its preferred alternative has enough pros, but at the same time it feels isolated with this opinion, because not enough



**Fig. 1** HUM-e process overview. *Being signalled to or inquired depends on dissonance reduction strategies of other agents. Within one tick, it may not take place at all or take place multiple times

alters in the ego network chose the same alternative. To resolve a social dilemma, the agent signals to one of the alters, trying to convince them to change their mind (see *Signalling* sub-model for details).

A non-social dilemma occurs in any other instance of cognitive dissonance, when the agent is convinced that the preferred alternative is popular enough among alters but, in its opinion, has significant non-social disadvantages. To resolve the non-social dilemma, the agent inquires with alters in its ego network to find more advantages, or to make the already existing ones seem more important (see *Inquiring* sub-model for details).

**Attitude formation**

When forming the attitude towards the preferred alternative, the agent revises its beliefs based on the collected information about the preferences and beliefs of contacted alters. Once all the information is up to date, the agent compares the expected satisfaction and dissonance level of the considered option with an alternative and selects the more preferred choice. Once the choice is made by all agents, the time progresses by 1 tick—an abstract unit which triggers information exchange. After each tick, fear is depreciated by 1.

## 3.2 Design Concepts

### 3.2.1 Basic Principles

The HUM-e emotive-socio-cognitive architecture is based on three major theoretical pillars.

- Emotional contagion theory;
- Social identity & identity fusion theories; and
- Cognitive consistency theories.[1]

Emotional contagion theory (ECT) was formulated and developed by Elaine Hatfield and colleagues [14, 15]. The main premise of the theory is that when individuals attend to others, they spontaneously mimic each other's emotional expressions (facial, vocal, postural, and instrumental), and therefore synchronize emotional states. How fluidly emotions move from one person to another depends on individual differences in susceptibility of individuals involved. Susceptibility to emotional contagion is the tendency to automatically mimic and synchronize with the expressions of others and, through afferent feedback from the facial and/or skeletal muscular activity, to experience or "catch" the others' emotions ([16], p. 149). HUM-e uses this assumption and focuses on the spread of the emotion of fear through social networks.

---

[1] Cognitive consistency theories as theoretical foundations of HUMAT were described at length elsewhere [19], therefore we will focus our attention on the theoretical foundations to additions to HUMAT.

Fear as an emotion is distinct from the dissatisfied motive of personal health (i.e., cognitively realising that an alternative might be harmful to one's health), even though the two are closely connected—an agent can only become scared about its own health if it believes that choosing an alternative will threaten its health and well-being.

The architecture also includes mechanisms informed by social identity theory (SIT) and identity fusion theory (IFT), both of which shed light on the ways in which a person's sense of identity can affect their motives. SIT hypothesizes that individuals from different social groups attempt to differentiate themselves from one another because of pressures to evaluate their own group positively through ingroup/out-group comparisons [17]. These value laden social differentiations can ratchet up tension between groups, which can impact people's motivation to protect their group. IFT postulates that motivation toward extreme behaviours is enhanced when a person's sense of their group becomes functionally equivalent to their sense of self [18]. In HUM-e, group affiliation influences the attitude agents take when communicating, eventually driving their opinions (1) closer towards a consensus, (2) farther apart in an act of repulsion, or (3) causing agents to ignore each other.

### 3.2.2 Emergence

Emergent outputs of the agent-based model are aggregates of the characteristics of individual agents (popularity of the choice alternatives, perceived evaluation of the choice alternatives, and average dissonance level of the agents).

### 3.2.3 Objectives

Every agent chooses the more preferred alternative. Knowledge of HUM-es is represented as cognitions—beliefs about how satisfying each alternative will be for the relevant motives of the individual [20]. Overall evaluation of each alternative is a cumulative evaluation of all motives (see Eq. 1). HUM-es only change their mind if they receive information significantly changing their knowledge about how (dis)satisfying alternative choices are. The chosen alternative maximizes the individual's overall motive evaluation and minimizes the level of experienced cognitive dissonance (see *Interaction* for details).

**Interaction**

The process of attitude formation is supplemented by information exchange between agents, which is implemented as a dissonance reduction strategy used by the agents and can take two forms: signalling and inquiring (see *Signalling* and *Inquiring* submodels). Communication flow is complete and one-directional. Complete means that information regarding all beliefs (pros and cons of both alternatives) is communicated in an interaction. Unidirectionality means that information is shared by one side (information source) and influences the other side of the conversation (target of the information). In signalling, the information flows from the signalling ego

(source) seeking to influence the alter (target). In inquiring, the alter (source) shares information with ego (the target). Each conversation between two agents will have one of three results: (1) bringing the target closer to a consensus with the source, (2) driving the repelled target farther away from the source's point of view or (3) the target ignoring the source (Table 1) depending on:

- Emotional state of fear of the source;
- Similarity of group affiliation between the source and the target;
- Similarity of the chosen alternative between the source and the target.

If the source of information is not afraid that the alternative threatens its health, and the source and target choose the same alternative, the target's new motive evaluation $\left( E_{a,j(t)}^{t_n+1} \right)$ will resemble the source's motive evaluation (consensus) to the extent the source is persuasive in the eyes of the target. Persuasiveness of the source depends on how close the beliefs of the two interlocutors originally are. The more similar views about how an alternative satisfies a motive the pair has the closer to 0.5 is the source's persuasiveness (weight of its opinion). Eventually, the target's new motive evaluation is a result of its previous views (minimum weight = 0.5), and the views of the information source (maximum weight = 0.5):

**Table 1** Results of conversations for the target's new attitude

| | | | g1 | | g2 | | g3 | |
|---|---|---|---|---|---|---|---|---|
| | | | a1 | a2 | a1 | a2 | a1 | a2 |
| *Source: no fear* | | | | | | | | |
| Target | g1 | a1 | Consensus | Consensus | Consensus | Ignore | Consensus | Repel |
| | | a2 | Consensus | Consensus | Ignore | Consensus | Repel | Consensus |
| | g2 | a1 | | | Consensus | Consensus | Consensus | Ignore |
| | | a2 | | | Consensus | Consensus | Ignore | Consensus |
| | g3 | a1 | | | | | Consensus | Consensus |
| | | a2 | | | | | Consensus | Consensus |
| *Source: in fear* | | | | | | | | |
| Target | g1 | a1 | Consensus + boost | Consensus + boost | Consensus + boost | Consensus | Consensus + boost | Consensus |
| | | a2 | Consensus + boost | Consensus + boost | Consensus | Consensus + boost | Consensus | Consensus + boost |
| | g2 | a1 | | | Consensus + boost | Consensus + boost | Consensus + boost | Consensus |
| | | a2 | | | Consensus + boost | Consensus + boost | Consensus | Consensus + boost |
| | g3 | a1 | | | | | Consensus + boost | Consensus + boost |
| | | a2 | | | | | Consensus + boost | Consensus + boost |

$$E_{a,j(t)}^{t_n+1} = 0.5 - \left(0.25\,abs\left(E_{a,j(t)}^{t_n} - E_{a,j(s)}^{t_n}\right) * E_{a,j(s)}^{t_n}\right)$$
$$+ 1 - \left(0.5 - \left(0.25\,abs\left(E_{a,j(t)}^{t_n} - E_{a,j(s)}^{t_n}\right)\right) * E_{a,j(t)}^{t_n}\right) \quad (2)$$

If the agents choose different alternatives, but the groups they affiliate with are not in conflict (i.e., they belong to $g_1$ and $g_2$ or to $g_2$ and $g_3$), the target simply ignores what the source is saying:

$$E_{a,j(t)}^{t_n+1} = E_{a,j(t)}^{t_n} \quad (3)$$

If they choose different alternatives and the source and the target identify with antagonistic groups (i.e., $g_1$ and $g_3$), the target will be repelled by what the source is saying, and will be driven away from source's beliefs to the same extent it would have been drawn towards the source's point of view, had the source belonged to the target's in-group:

$$E_{a,j(t)}^{t_n+1} = 2 * E_{a,j(t)}^{t_n} - \left(0.5 - \left(0.25\,abs\left(E_{a,j(t)}^{t_n} - E_{a,j(s)}^{t_n}\right) * E_{a,j(s)}^{t_n}\right)\right)$$
$$+ 1 - \left(0.5 - \left(0.25\,abs\left(E_{a,j(t)}^{t_n} - E_{a,j(s)}^{t_n}\right)\right) * E_{a,j(t)}^{t_n}\right)\right) \quad (4)$$

If the source is afraid that a choice alternative threatens its health and well-being, the genuine expression of emotions will make it seem convincing to the target. Effectively, the conversation will be driven by the goal of consensus (see Eq. 1), irrespectively of group affiliation and chosen alternative. If the agents choose the same alternative and/or perceive each other as in-group members, the scared source will get an extra persuasiveness boost—a maximum of 0.05, depending on the level of target's fear contagion $\left(C_{j(t)}\right)$. Being in fear for your life is therefore assumed to cross barriers built by group affiliations:

$$E_{a,j(t)}^{t_n+1} = \left(0.5 - \left(0.25\,abs\left(E_{a,j(t)}^{t_n} - E_{a,j(s)}^{t_n}\right) + \left(0.05 * C_{j(t)}\right)\right) * E_{a,j(s)}^{t_n}\right)$$
$$+ 1 - \left(0.5 - \left(0.25\,abs\left(E_{a,j(t)}^{t_n} - E_{a,j(s)}^{t_n}\right) + \left(0.05 * C_{j(t)}\right)\right) * E_{a,j(t)}^{t_n}\right)$$
$$(5)$$

### 3.3  Details

**Initialization**

The HUM-e architecture can be calibrated to represent human decision-making in a wide range of cases, where beliefs and fear for one's health/safety/well-being spread

**Fig. 2** Initialization of agents

through communities. Initiation of the "empty" architecture draws attention to the dimensions of calibration and signal data needs.

**World initialization**

At initiation, a selected number of agents are distributed over a random radius from the centre of the world. The agents link with a maximum of 20% other agents located within a distance of 20 patches, so that each agent has at least one alter in the ego network. For calibration to a specific case, GIS data and relevant networking can be implemented.

**HUM-e initialization**

Once the world setup is finalized, the agents are initialized (Fig. 2).

**Non-social motives**: The modeller chooses the properties of the distribution of the initial non-social motive evaluations (experiential motive ($e_1$), values ($v_1$) and health motive ($h_1$)) for each choice alternative. In the non-calibrated architecture, this can be achieved in two ways: either by choosing from a list of possible choices (heavily right-skewed, slightly right skewed, normal, slightly left skewed, heavily left skewed), or by manually inputting the parameters of the beta distribution in the code. Quantitative empirical data can be used to calibrate agent's initial evaluations. In the absence of quantitative data, qualitative insight can also be utilized.

Fear is a temporary state grounded in cognitive beliefs, which lasts for 3 consecutive ticks once it's triggered. All agents who find at least one choice alternative scary with respect to their health $\left( E_{a,h_1,j}^{t_n} \leq -0.7 \right)$ set their fear to a random integer between 1 and 3 (to avoid all agents synchronizing on three ticks at initiation). Other individual differences initialized at the setup stage follow a random normal distribution ($\mu = 0.5$; $\sigma = 0.14$) ranging from 0 to 1:

- fear contagion $\left( C_{j(t)} \right)$—the individual difference governing how susceptible to fear the target of the information is. The most fear susceptible agents (fear contagion = 1) perceive the source of information as more persuasive (to a maximum boost of 5%) when the source communicates about an alternative it's afraid of.
- social tolerance—the individual difference governing the need to belong to a group of similarly choosing alters. The most socially needy agents (social need

importance $= 1$) require all alters in the ego network to choose the same alternative as they do. On the other end of the spectrum, agents with low social tolerance are individualistically-minded and do not give into social pressure easily. If the fraction of alters who choose the same alternative as ego is smaller than the social need importance, the agent's social need is evaluated negatively, and the agent may engage in signalling to the most gullible alter to change its attitude.

- dissonance tolerance ($T_j$)—the individual difference governing how much cognitive dissonance the agent is comfortable with without the need to implement dissonance reduction strategies (signalling and inquiring).

**Basic choice**: Agents make their basic choice between alternatives. The chosen alternative is cumulatively evaluated highest with respect to the evaluations of initialized non-social motives.

**Initial choice**: Agents set up representations of alters in their ego networks. For each alter, the ego stores the information about the identified ([who] of alter), the fact of having had inquired of them already (0 at initiation), the fact of having signalled to the alter already (0 at initiation), perception of alter choice (80% correct—corrected when information exchange takes place), the distance between ego and alter group affiliations (perfect knowledge), whether alter is afraid of alternative $a_1$(0 at initiation—changed when information exchange takes place), and whether alter is afraid of alternative $a_2$(0 at initiation—changed when information exchange takes place). For calibration purposes, the list can be extended. The ego counts the alters, who, in its perception, chose the same alternative. This is the basis for the evaluation of the social motive. If the fraction of alters who choose the same alternative exceeds ego's social tolerance, ego evaluates the alternative positively. Similar to evaluations of other motives, social evaluation is normalized between $-1$ and 1 with social tolerance level $= 0$.

Afterwards, the agent calculates cognitive dissonances and identifies choice dilemmas. Agents experience cognitive dissonance when a choice alternative is perceived as internally inconsistent—it has both pros (positively evaluates some motives) and cons (negatively evaluates other motives). The amount of dissonance a choice alternative evokes:

$$D_{a,j}^{t_n} = 2d/d + c$$

where:

$d$—dissonant cognitions $\left(min\left(\sum E_{a,j}^{t_n} > 0,\ \sum E_{a,j}^{t_n} < 0\right)\right)$,

$c$—consonant cognitions $\left(max\left(\sum E_{a,j}^{t_n} > 0,\ \sum E_{a,j}^{t_n} < 0\right)\right)$.

Most often dissonant cognitions are suppressed or ignored as an effective dissonance reduction strategy [21]. Therefore, in HUM-e, dissonance needs to be actively resolved when it exceeds the individual's tolerance threshold ($T_j$). As a consequence of unignorable dissonance, the agents either signal or inquire—depending on the type of dilemma they face. A social dilemma occurs when the social evaluation of a choice

alternative is negative, and evaluation of at least one other motive is positive (Table 2). Non-social dilemma occurs when social evaluation is positive and at least one other motive is evaluated negatively.

Finally, the agent makes an initial choice, which is based on all motive evaluations. If the evaluations of both alternatives are sufficiently similar with respect to cumulative evaluation, the agent chooses the alternative which is not scary $\left( E_{a,h_1,j}^{t_n} > 0 \right)$. If both alternatives are similar with respect to how scary they are, the agent chooses the alternative which is less dissonant. If both alternatives are similarly dissonant, the agent sticks to the alternative chosen previously. Cognitive dissonance is a motivational force for a change in knowledge [22] or behaviour [23]. Occurrence of cognitive dissonance leads to a psychologically unpleasant state of facing a dilemma.

**Alter representations**: The agents update alter representations (in the same procedure implemented in the *Initial choice*) and set inquiring and signalling lists (see *Signalling* and *Inquiring* sub-models for details).

**Dissonances**: Based on new information from social networks, agents update their cognitive dissonances and dilemmas, and are ready to implement dissonance reduction strategies once the model starts running (See *Process overview and scheduling* for details).

**Table 2** Social and non-social dilemmas in HUM-e

| Experiential motive $e_1$ | Values $v_1$ | Health motive $h_1$ | Social motive $s_1$ | Dilemma |
|---|---|---|---|---|
| $\geq 0$ | $\geq 0$ | $\geq 0$ | $< 0$ | Social dilemma |
| $\geq 0$ | $\geq 0$ | $< 0$ | $< 0$ | Social dilemma |
| $\geq 0$ | $< 0$ | $\geq 0$ | $< 0$ | Social dilemma |
| $< 0$ | $\geq 0$ | $\geq 0$ | $< 0$ | Social dilemma |
| $\geq 0$ | $< 0$ | $< 0$ | $< 0$ | Social dilemma |
| $< 0$ | $< 0$ | $\geq 0$ | $< 0$ | Social dilemma |
| $< 0$ | $< 0$ | $< 0$ | $< 0$ | No dilemma |
| $< 0$ | $< 0$ | $< 0$ | $\geq 0$ | Non-social dilemma |
| $< 0$ | $< 0$ | $\geq 0$ | $\geq 0$ | Non-social dilemma |
| $< 0$ | $\geq 0$ | $< 0$ | $\geq 0$ | Non-social dilemma |
| $\geq 0$ | $< 0$ | $< 0$ | $\geq 0$ | Non-social dilemma |
| $< 0$ | $\geq 0$ | $\geq 0$ | $\geq 0$ | Non-social dilemma |
| $\geq 0$ | $\geq 0$ | $< 0$ | $\geq 0$ | Non-social dilemma |
| $\geq 0$ | $\geq 0$ | $\geq 0$ | $\geq 0$ | No dilemma |

**Sub-models**

**Signalling**

If the preferred alternative is not popular enough among alters $\left( E_{s_1, j}^{t_n} < 0 \right)$ and the dissonance level exceeds the agent's tolerance threshold, ego signals to its most gullible alter with an opposite preference and tries to convince it to change its mind. To this end, ego copies the alter representation list and stores it as the signalling list after applying appropriate sorting and identifying in the beginning of the list the alters who: (1) have not yet been signalled, (2) choose a different alternative, and (3) have the closest group affiliation. Once the most gullible agent is identified as first on the signalling list, the conversation between agents starts and the target of the information listens to the signalling source, changes its beliefs about how different alternatives satisfy individual motives (see *Interaction* for details), and forms a new attitude (see *Attitude formation* for details). Subsequently, the signalling agent updates the information about the alter and forms a new attitude. As a result of the conversation, both agents have updated information about their interlocutor (updated choice alternative, the facts of having signalled to and being signalled to) and formed new attitudes (including recalculating dissonances, identifying dilemmas, and possibly making a new choice).

**Inquiring**

If the preferred alternative is popular enough among alters, but evokes unignorable dissonance for other reasons, the ego will actively look for advice in its social network. To this end, ego copies the alter representation list and stores it as the inquiring list after applying appropriate sorting and identifying in the beginning of the list the alters who: (1) have not yet been objects of inquiry, (2) choose the same alternative, and (3) have closest group affiliations. Ego perceives such an alter as the most persuasive and becomes a target of the information sharing (see *Interaction* for details). As a result of the interaction, both agents form new attitudes (see *Attitude formation* for details).

# References

1. Grimm, V., Berger, U., DeAngelis, D., Polhill, J., Giske, J., Railsback, S.: The ODD protocol: a review and first update. Ecol. Model. **221**(23), 2760–2768 (2010). https://doi.org/10.1016/j.ecolmodel.2010.08.019S
2. Simon, H.: Rational decision making in business organizations. Am. Econ. Rev. **69**(4), 493–513 (1979). http://www.jstor.org/stable/1808698
3. Merkuur, R.: Simulating human routines. Integrating social practice theory in agent-based models (2021). https://doi.org/10.4233/uuid:4b70aa0a-8c13-421d-9043-6274311df2aa
4. Gratch, J., Marsella, S.: Appraisal models. In: Calvo, R., D'Mello, S., Gratch, J., Kappas, A. (eds.) The Oxford Handbook of Affective Computing, 54–67 (2014)
5. Hudlicka, E.: Two sides of appraisal: Implementing appraisal and its consequences within a cognitive architecture. Papers from the 2004 AAAI Spring Symposium (2004).

6. Schröder, T., Hoey, J., Rogers, K.B.: Modeling dynamic identities and uncertainty in social interactions: Bayesian affect control theory. Am. Sociol. Rev. **81**(4), 828–855 (2016)
7. Dancy, C.L.: ACT-RΦ: a cognitive architecture with physiology and affect. Biol. Inspired Cogn. Archit. **6**, 40–45 (2013)
8. Allen, J., Sun, R.: Emotion contagion in a cognitive architecture. In: 2016 IEEE Symposium Series on Computational Intelligence (SSCI), pp. 1–8 (2016)
9. Jiang, H., Vidal, J.M., Huhns, M.N.: EBDI: an architecture for emotional agents. In: Proceedings of the 6th International Joint Conference on Autonomous Agents and Multiagent Systems, pp. 1–3 (2007)
10. Puica, M.-A., Florea, A.-M.: Emotional belief-desire-intention agent model: previous work and proposed architecture. Int. J. Adv. Res. Artif. Intell. **2**(2), 1–8 (2013)
11. Zickfeld, J.H., Schubert, T.W., Herting, A.K., Grahe, J., Faasse, K.: Correlates of health-protective behavior during the initial days of the COVID-19 outbreak in Norway. Front. Psychol. **11**(564083) (2020). https://doi.org/10.3389/fpsyg.2020.564083
12. Harper, C.A., Satchell, L., Fido, D., Latzman, R.: Functional fear predicts public health compliance in the COVID-19 pandemic. Int. J. Ment. Health Addict. **27**, 1–14 (2020). https://doi.org/10.1017/dmp.2020.338
13. Wheaton, M.G., Prikhidko, A., Messner, G.R.: Is fear of COVID-19 contagious? The effects of emotion contagion and social media use on anxiety in response to the coronavirus pandemic. Front. Psychol. **11**(567379) (2021). https://doi.org/10.3389/fpsyg.2020.567379
14. Hatfield, E., Cacioppo, J.T., Rapson, R.L.: Emotional contagion. Curr. Dir. Psychol. Sci. **2**(3), 96–100 (1993). https://doi.org/10.1111/1467-8721.ep10770953
15. Hatfield, E., Bensman, L., Thornton, P.D., Rapson, R.L.: New perspectives on emotional contagion: a review of classic and recent research on facial mimicry and contagion. Interpersona **8**, 159–179 (2014)
16. Doherty, R.W.: The emotional contagion scale: a measure of individual differences. J. Nonverbal Behav. **21**(2), 131–154 (1997)
17. Tajfel, H. (ed.): Social Identity and Intergroup Relations, reissue edn. Cambridge University Press (2010)
18. Swann, W.B., Buhrmester, M.D.: Identity fusion. Curr. Dir. Psychol. Sci. **24**(1), 52–57 (2015). https://doi.org/10.1177/0963721414551363
19. Antosz, P., Jager, W., Polhill, G.: Simulation model implementing different relevant layers of social innovation, human choice behaviour and habitual structures. SMARTEES Deliverable (2019)
20. Festinger, L.: Reflections on cognitive dissonance: 30 years later. In: Harmon-Jones, E., Mills, J. (eds.) Cognitive Dissonance: Progress on a pivotal theory in social psychology. American Psychological Association, Washington (1999)
21. McGrath, A.: Dealing with dissonance: a review of cognitive dissonance reduction. Soc. Pers. Psychol. Compass **11**(12), 1–17 (2017)
22. Festinger, L.: A theory of social comparison processes. Human Relat. **7**, 117–140 (1954)
23. Harmon-Jones, E., Harmon-Jones, C.: Testing the action-based model of cognitive dissonance: the effect of action orientation on postdecisional attitudes. Pers. Soc. Psychol. Bull. **28**(6), 711–723 (2002)

# Identity Drives Polarization: Advancing the Hegselmann-Krause Model by Identity Groups

**František Kalvas** ⓘ**, Ashwin Ramaswamy** ⓘ**, and Michael D. Slater** ⓘ

**Abstract** In this article we describe an Agent-Based Model that extends the Hegselmann-Krause model of opinion dynamics to study the role of social identity in opinion polarization. In our model, an agent's social identity is a function of two things—the agent's opinion in relation to those of the other agents, and the observer's sensitivity to the tightness of clustering. We implement this by first selecting a subset of the agent population that are deemed to have close neighbors, and then using Louvain community detection to find identity groups. At every time step, agents only consider the opinions of other agents within their identity group that also fall within their Hegselmann-Krause opinion boundary, ε. We show that our dynamic implementation of social identity systematically modulates the relationship between average ε and polarization.

**Keywords** Bounded confidence model · Dynamic identity · Polarization

## 1 Introduction

The process of consensus formation in public opinion is at least partially believed to be impacted by social interactions. Individuals gather information about the world, other individuals, and societal structures through conversations with each other. They also learn about accepted norms and normative evaluations of individuals and situations through interactions with others. Through this process they ultimately form

F. Kalvas (✉)
University of West Bohemia, Univerzitní 8, 301 00 Plzeň, Czech Republic
e-mail: kalvas@kss.zcu.cz

A. Ramaswamy
Pratiksha Nagar H-16, Mumbai 400022, India

M. D. Slater
The Ohio State University, 154 North Oval Mall, Columbus, OH 43210-1339, USA

their own beliefs and opinions about relevant issues using social information as one of the inputs. People with differing points of view may reconcile their differences through conversation by either learning to adopt the other's views, or convincing the other of one's own or by resisting opinion change.

The Hegselmann-Krause (HK) model of opinion dynamics is a bounded-confidence model with continuous real-valued opinions [1]. Classically, the HK imposes a constraint that allows a listening agent to only consider other agents whose opinion falls within a distance of a boundary parameter (commonly denoted by ε) of the listening agent. The listening agent updates its opinion to the average value of the opinions of all such agents. Different system parameters and initial conditions can cause the HK system to produce consensus, polarization, or fractured states.

A number of theoretical properties of the HK model have been studied such as the probability [2] and kinetics of consensus [3], the role of noise [4, 5], heterogeneity in ε [6], adding more dimensions to the opinion space [3, 7], and how the presence of social network constraints influences dynamics [8]. Some other studies have extended the HK model by adding to the dynamics new features such as the presence of agenda-setting 'leaders' [9] or extremists [10].

We took a slightly different approach to advancing the HK model, by introducing an additional component to the dynamics that simulates the role of social identity groups on the asymptotic behavior of the system. Social Identity Theory (SIT) proposes that pairwise inter-personal interactions are relevant but insufficient to explaining the collective dynamics of a society, and that perceived group identities influence one's behavior towards another [11]. Identities may help individuals understand and approximate a complex landscape of public opinion and interests by reducing nuances and variances into simplified labels. We aimed to study the relevance of social identities to polarization—a qualitative state of the system where the opinions of all agents tend to be split into two antagonistic camps—given its sociological significance as a commonly occurring state of public opinion [12, 13].

Both assumed (by the self) and perceived (of others) identities are known to influence one's opinion. For example, a study by Wojcieszak and Garrett found that priming national identity, and exposure to anti-immigration news increases reported anti-immigration sentiment among anti-immigration participants [14]. We follow the Reinforcing Spirals Model [15, 16] in proposing that salience of social identity and the degree to which there exists closed vs open communication norms are major drivers of polarization in a dynamic model.

Consistent with SIT, we treat opinions and identities as interacting components of social behavior that are both relevant for dynamics. Therefore, we model the formation of social identities as an emergent process in the opinion space. Agents look at the entire opinion space to find groups of agents that are well-clustered, and assign identities to these clusters. Then they update their opinions using only the inputs from agents that are both within their own identity group, and also satisfy the HK opinion boundary. Our model thus assumes that social identity acts as an additional filter for agents as they select other agents to seek consensus with at each step. Therefore, an agent might ignore another's opinion either because their opinions are too far from each other, or because they perceive the other agent to be

in a different identity group. Importantly, identities are not pre-assigned to agents—rather they are inferred from the opinion positions of the entire population. Some agents may "see" different identity groups than others. Moreover, identity groups might evolve as opinions of agents evolve—as agents move through the opinion space the groups might merge, shift and break up.

We needed a plausible algorithm to dynamically assign agents to identity groups based on their opinion positions in the opinion space. For this, we needed to consider what conditions must be satisfied for agents to be said to form an identity group based on opinions. Firstly, for an identity group to be said to exist, there must be at least a few agents showing a high degree of proximity to one another in the opinion space. Secondly, for an agent to be considered as part of an identity group, she must demonstrate sufficient similarity to the identity group's ideology. Thirdly, the identity group must not only be defined by the proximity of the opinions of its own members, but must also be sufficiently far from agents it excludes. In other words an identity group isn't defined just by the oneness of its members, but must also take into account the otherness of agents it excludes.

Our algorithm for identity group detection follows a similar logic as detailed above. Identity groups are detected in the opinion space by considering only those agents that have enough sufficiently like-minded agents in the opinion space. In this subset of non-isolated agents, the detector applies a Louvain Community Detection (LCD) [17] algorithm, which is our implementation of a general mechanism that lets agents automatically detect the existence of identity groups from information about the spread of opinions in the population.

An important parameter in the process above controls what we mean by "sufficient like-mindedness" in the filtering step. We call this parameter "Salience of Proximity in Identity-Relevant Opinions" (SPIRO), since it defines which pairs of agents are close enough in opinion space to be relevant for LCD, and treat it as an experimental variable. SPIRO is a property of the detector—as agents look around in the opinion space and detect identity groups, they may be differentially sensitive to agents clustering close together.

In this article we present five hierarchically related models, of which the last two include social identity effects. We do this to introduce not only our implementation of social identity, but also other model features and variables we believe may have interesting effects alongside identity. In Sect. 2 we discuss our methods, including their components (Sect. 2.1), the model variants (Sect. 2.2), and our variables of interest (Sect. 2.3). In Sect. 3 we present evidence that the presence of identity drives polarization, along with some preliminary results involving other variables. In Sect. 4 we interpret these data and present plans for future work with these models.

## 2 Methods (Code and Data are Available [18])

### 2.1 Model Components

**Hegselmann-Krause Dynamics with conformity.** An agent's opinion is represented as a real number between $-1$ and $+1$, implemented at the resolution of 3 decimal places. Opinions of all agents are updated at every time-step based on their previous opinion and the opinions of influencers according to the rule:

$$o_i(t) = o_i(t-1) + \alpha_i \left[ \frac{1}{|N_i(t)|} \sum_{j \in N_i(t)} o_j(t-1) - o_i(t-1) \right] \tag{1}$$

where,

$o_i(t) \in [-1, +1]$ is the opinion of agent $i$ at time $t$.

$\alpha_i \in [0, 1]$ is the conformity parameter, it controls how quickly agent $i$ moves towards the found consensus.

$N_i(t)$ is the neighborhood of agent with index $i$ at time $t$.

$$N_i(t) = \{j : |o_j(t-1) - o_i(t-1)| \leq \varepsilon_i\} \tag{2}$$

$\varepsilon_i \in [0, 1]$ is the boundary parameter and tells us the maximum dissimilarity in opinion agent $i$ can accommodate. Note that $\varepsilon_i$ is normalized—it is measured as a fraction of the maximum possible distance, i.e. $\varepsilon_i = 1$ means that agent $i$ with opinion $o_i = -1$ also takes into account agent $j$ with opinion $o_j = +1$.

Thus $N_i(t)$ is the set of all agents (including the listening agent itself) whose opinion fall within a distance of the boundary parameter $\varepsilon_i$ of the listening agent.

**Social Identity Boundary.** In the model with social identity, an agent only listens to another agent who additionally also shares the same identity group as oneself at each time step.

Let $Id_i(t)$ represent the index of the identity group of agent $i$ at time $t$. Therefore, the neighborhood of an agent $N_i(t)$ is redefined as:

$$N_i(t) = \{j : |o_j(t-1) - o_i(t-1)| < \varepsilon_i\} \cap \{j : Id_j(t) = Id_i(t)\} \tag{3}$$

**Identity Group Assignment.** The identity groups are dynamically updated at every time step as follows:

Firstly, we convert the opinion space into an equivalent weighted full network (the 'Proximity Network') by representing each agent by a node creating a weighted link between every pair of agents. The weight of each link is given by:

$$w_{i,j}(t) = 1 - d(o_i(t), o_j(t)) \tag{4}$$

where,

$w_{i,j}(t) \in [0, 1]$ is the weight of the link between nodes $i$ and $j$ at time $t$.

$d(a, b)$ is the Euclidean distance between points a and b in the opinion space, normalized by the maximum theoretical distance in the opinion space. Therefore,

$$d\big(o_i(t),\ o_j(t)\big) = \frac{\big|o_i(t) - o_j(t)\big|}{2} \tag{5}$$

Thus a weight of 1 means the two linked agents have identical opinions, while a weight of 0 means they are maximally dissimilar.

We then perform community detection on a subset of the Proximity Network, keeping only edges of sufficient weight and nodes sufficiently connected by such edges. We use a SPIRO-thresholded definition of which edges' weights are sufficient, and we keep only nodes connected by 2 or more such edges, along with only edges of sufficient weight to these nodes. We then perform LCD on this sub-graph. In practice, "sufficiently connected" edges are edges whose weight in the Proximity Network equals or exceeds the perceiving agent's SPIRO value. Thus, higher SPIRO values would mean we tend to return fewer nodes and links after these reduction steps.

In order to ensure every agent is assigned to an identity group, we follow up LCD with k-means clustering as follows—we consider the number of detected communities after SPIRO-thresholding and LCD on the Proximity Network, and compute the opinion centroid of the set of agents corresponding to each community. We use the number of communities and the centroids thus found as initial values to the k-means clustering algorithm which is performed on the entire agent population (including those excluded before LCD). Every excluded agent is initially assigned to the cluster whose centroid is closest to it. k-means clustering is repeated on the opinion space thereafter until the centroids converge. Thus, every agent is assigned to an identity group.

**Global versus individual detection of identity groups**. We wanted to simulate the possibility of different agents being differently sensitive to identity-related information from the opinion space—in our model this translates to agents having different SPIRO values (see Sect. 2.2, model VBVI). Implementing this directly would mean running the Louvain algorithm several times at every time step, making the simulation computationally very expensive. To make the process more efficient, we segmented the agent population into eight partitions, each having its own pre-defined SPIRO value. Although the number and index of the agents assigned to each partition may vary across simulations, every partition—and therefore every agent—can take on SPIRO values only from the set: {0.15, 0.25, 0.35, 0.45, 0.55, 0.65, 0.75, 0.85}.

To determine which agent gets assigned to which SPIRO value, we implemented an approximation of a discrete normal SPIRO distribution as follows: During simulation set up, every agent samples a value $x_i$ from a normal distribution with mean

$\mu_{SPIRO}$ and standard deviation $\sigma_{SPIRO}$. If $x_i \notin [0, 1]$, its sampling is repeated until $x_i \in [0, 1]$. The SPIRO of the agent $i$ is given by the closest possible value to $x_i$ from the set of valid SPIRO's given above.

## 2.2 Model Variants

We ran 2,504,964 simulations in total spanning 5 variants of the HK model. The models are described:

*Deterministic Start HK Model (DHK):* Initial opinions of agents are a set of evenly spaced real numbers between $[-1, +1]$. Agents have the same confidence boundary $\varepsilon$.

*Randomized Start HK Model (RHK):* Initial opinions of agents are uniformly distributed real values in the interval $[-1, +1]$. Agents have the same confidence boundary $\varepsilon$.

*Heterogeneous Boundary Model (VB):* Agents have individualized confidence boundaries and conformities. The confidence boundary $\varepsilon_i$ of an agent is obtained from a truncated normal distribution as follows: Every agent samples a value $\varepsilon_i$ from a normal distribution with mean $\mu_\varepsilon$ and standard deviation $\sigma_\varepsilon$. If $\varepsilon_i \notin [0, 1]$, its sampling is repeated until $\varepsilon_i \in [0, 1]$. $\alpha_i$ is also sampled with an identical method as $\varepsilon_i$, with mean $\mu_\alpha$ and standard deviation $\sigma_\alpha$.

*Heterogeneous Boundary with Identity (VBI):* Agents only communicate within their identity groups, which are assigned at the beginning at every time-step via a common identity group assignment step as outlined in Sect. 2.1. This assignment is parametrized by the common SPIRO value, which determines the tightness of identity groups thus formed.

*Heterogeneous Boundary with Heterogeneous Identity (VBVI):* Agents only communicate within their perceived identity groups, but they may be inconsistent across agents. This is done by relaxing the assumption of a single SPIRO value for the entire population as follows:

1. At the beginning of the simulation all agents are assigned an individualized SPIRO value as described in Sect. 2.1. This is done to allow for heterogeneous identity effects while keeping the model computationally efficient.
2. At the beginning of every time step, one instance of the identity group assignment step outlined in Sect. 2.1 is run for each partition.
3. The detected identity groups for each partition are then inherited by each agent within the partition. Thus, all the agents in a partition perceive a common set of identity groups.

The above models are hierarchically related, in that every subsequent model in the list above inherits features of the previous models (exception: RHK does not inherit

the regularly-spaced initial opinion space condition from DHK). Therefore, VBI and VBVI both have normally distributed $\varepsilon$ values for instance. We ran each simulation for 365 time steps, or until consensus is reached, whichever is earlier.

## 2.3 Variables

**Independent Variables**. Besides $\mu_\varepsilon$, $\sigma_\varepsilon$, $\mu_\alpha$, $\sigma_\alpha$, $\mu_{SPIRO}$, and $\sigma_{SPIRO}$ which are defined in Sects. 2.1 and 2.2, we also included the following two variables in our experimental design since we were also interested in studying some robustness properties of the HK model for a related study:

*Evenness or Oddness of population size:* Population size is either N = 100 or N = 101.

*Randomness of initial opinion distribution:* The initial opinion of agents is either drawn uniformly at random (Random_start? = TRUE), or can assume equally spaced out values in the interval of $[-1, +1]$ (Random_start? = FALSE).

Note: in models with no variability of some parameter p, $\mu_p$ stands in for the common value of p.

**Dependent Measure—Polarization**. To measure polarization we adapt the Equal Size Binary Grouping (ESBG) algorithm from Tang et al. [19], which gives a continuous-valued metric we call ESBG Polarization, or just ESBG. The ESBG measure is based on the ideal type of maximally polarized community. Such a community is divided in two camps of equal size. These camps are very homogenous, i.e. opinions of camp's members are the same, but these camps are on opposite poles of opinion scale, i.e. the distance of camps in opinion space is maximal. To reflect this ideal type, ESBG firstly divides the population in two groups by a specific version of k-means clustering algorithm. This algorithm divides the population in two groups of equal size, but on the other hand it minimizes opinion heterogeneity of these forcibly created groups. Then ESBG computes distance of group centroids and mean deviation of groups' members' opinions around respective centroids. Then ESBG value is computed as centroids' distance divided by sum of 1 and mean deviations of both clusters. Centroids' distance and mean deviations of both clusters are normalized by maximum possible distance which ensures that the resulting ESBG is between 0 and 1, where 0 signifies perfect consensus and 1 signifies complete polarization.

$$ESBG = \frac{Norm(B)}{1 + Norm(w_1) + Norm(w_2)} \tag{6}$$

where,

$$Norm(x) = \frac{x}{\sqrt{4 \times Number\ of\ Opinion\ Dimensions}} = \frac{x}{2} \qquad (7)$$

$B$ = Absolute distance between the two cluster centroids

$w_i$ = Total mean deviation of agent opinions of cluster $i$ from its centroid

**Analysis**: We performed multiple regression for our dependent measure on the experimental variables of interest: $\mu_\varepsilon$, $\sigma_\varepsilon$, $\mu_\alpha$, $\sigma_\alpha$, $\mu_{SPIRO}$, $\sigma_{SPIRO}$, Evenness of Population Size, and Randomness of initial opinions. To avoid making assumptions about linearity of relationships we treated each variable as a factor. In our results section we report mean ESBG value of all simulations run for a given combination of variables.

## 3   Results

The relationships between polarization and $\mu_\varepsilon$ of each of our models show qualitative differences (Fig. 1). Firstly, we observe that the two models with dynamically updated identity groups (VBI and VBVI) maintain polarized states for much higher values of $\mu_\varepsilon$ than the other models. Secondly, we observe the lowest polarization in the Heterogeneous Boundary Model (VB), the difference in polarization is striking and significant especially for lower values of $\mu_\varepsilon$ (approximately in interval 0.10–0.23). Thirdly, we observe the effect of deterministic starting conditions: polarization produced by the DHK model in response to $\mu_\varepsilon$ values qualitatively dramatically differs from all other models based on or employing random start conditions. Fourthly, we observe that $\sigma_{SPIRO}$ has a negligible effect—models VBI and VBVI differ just slightly and they reach maximal difference only for the highest investigated value of $\mu_\varepsilon$. Fifthly, we observe that $\mu_\alpha$ and Size of Population (N) have effect on models not employing heterogenous Boundary (DHK and RHK), models VB, VBI and VBVI seem to qualitatively keep their behavior despite the values of $\mu_\alpha$ and N.

Here we report that evenness appears to drive a qualitative change in the Polarization-Boundary relationship only when the initial condition is not randomized (DHK). We originally investigated the effect of population size. We surprisingly found that size itself does not matter much, but what matters for DHK was whether the population size is even or odd. For example, even for DHK it had almost no effect whether the size of population was 20, 100, or 256 agents, but it had a substantive effect whether the size was 21 instead of 20, or 101 instead of 100, or 257 instead of 256 agents. For the final presentation of our analyses in this paper we chose N = {100, 101}, since these sizes spot the effect of evenness and are heavily used in the canon of literature. We intend to explore this methodological issue further in a subsequent paper (in preparation).

**Fig. 1** ESBG-$\mu_\varepsilon$ relationship for each model. Panels represent different conditions of population evenness and conformity. Ordinate in each panel is the Mean ESBG at the end of all simulations with the given parameter combination

Mean ESBG polarization differs across our models in the following way: VB < DHK < RHK < VBVI < VBI (Table 1). The two models with identity in them have the highest mean polarization—showing that identity drives polarization and impedes consensus. In all the models, $\mu_\varepsilon$ is negatively associated with polarization as expected (Table 2).

A consistent finding throughout our analyses is that higher $\sigma_\varepsilon$ brings down polarization dramatically (Table 2), and its influence is stronger than that of the mean boundary. This is also evident in Figs. 2 and 3. We interpret this as an unbalanced mitigating influence of agents with higher-than-average boundaries (see discussion). $\sigma_{SPIRO}$ also lowers polarization, although far not as strongly as $\sigma_\varepsilon$.

The main drivers of polarization are $\mu_{SPIRO}$, $\sigma_\varepsilon$, and $\mu_\varepsilon$. This can also be seen in Figs. 2 and 3 for a model with heterogeneous identity (VBVI). Interestingly however, $\mu_{SPIRO}$ systematically modulates the relationship between $\mu_\varepsilon$ and polarization (Fig. 3). For $\mu_{SPIRO}$ values from 0.25 to 0.61, $\mu_{SPIRO}$ is positively associated

**Table 1** Summary statistics for ESBG in different models

| Model | N | Min | Max | IQR | Median | Mean | SD | SE | CI |
|---|---|---|---|---|---|---|---|---|---|
| DHK | 84 | 0 | 0.419 | 0.361 | 0.282 | 0.199 | 0.177 | 0.019 | 0.038 |
| RHK | 5040 | 0 | 0.534 | 0.371 | 0.304 | 0.242 | 0.167 | 0.002 | 0.005 |
| VB | 80,640 | 0 | 0.872 | 0.251 | 0.026 | 0.114 | 0.157 | 0.001 | 0.001 |
| VBI | 483,840 | 0 | 0.937 | 0.154 | 0.408 | 0.378 | 0.177 | 0.000 | 0.000 |
| VBVI | 1,935,360 | 0 | 0.940 | 0.208 | 0.405 | 0.354 | 0.195 | 0.000 | 0.000 |

**Table 2** Regression on ESBG in model VBVI. (N = 460,800)

|                          | Estimate  | Std. Error | t value    | Pr(>|t|) |
|--------------------------|-----------|------------|------------|----------|
| Intercept                | 0.401     | 0.001      | 354.267    | 0.000    |
| $\sigma_{SPIRO}$ *(contrast: 0)* |   |            |            |          |
| 0.05                     | − 0.015   | 0.001      | − 21.450   | 0.000    |
| 0.10                     | − 0.027   | 0.001      | − 38.374   | 0.000    |
| 0.15                     | − 0.036   | 0.001      | − 50.784   | 0.000    |
| $\mu_{SPIRO}$ *(contrast: 0.25)* |  |            |            |          |
| 0.37                     | 0.057     | 0.001      | 65.217     | 0.000    |
| 0.49                     | 0.110     | 0.001      | 125.936    | 0.000    |
| 0.61                     | 0.150     | 0.001      | 171.538    | 0.000    |
| 0.73                     | 0.076     | 0.001      | 86.889     | 0.000    |
| 0.85                     | 0.091     | 0.001      | 103.378    | 0.000    |
| $\sigma_{\varepsilon}$ *(contrast: 0)* | |        |            |          |
| 0.05                     | − 0.021   | 0.001      | − 28.738   | 0.000    |
| 0.10                     | − 0.129   | 0.001      | − 180.832  | 0.000    |
| 0.15                     | − 0.151   | 0.001      | − 210.761  | 0.000    |
| $\mu_{\varepsilon}$ *(contrast: 0.10)* | |        |            |          |
| 0.15                     | − 0.004   | 0.001      | − 5.087    | 0.000    |
| 0.20                     | − 0.029   | 0.001      | − 36.039   | 0.000    |
| 0.25                     | − 0.063   | 0.001      | − 78.256   | 0.000    |
| 0.30                     | − 0.112   | 0.001      | − 139.412  | 0.000    |
| *Random_start? (contrast: TRUE)* | |          |            |          |
| FALSE                    | 0.022     | 0.001      | 43.744     | 0.000    |
| $\sigma_{\alpha}$ *(contrast: 0)* | |          |            |          |
| 0.10                     | − 0.001   | 0.001      | − 1.230    | 0.219    |
| $\mu_{\alpha}$ *(contrast: 0.20)* | |          |            |          |
| 0.80                     | − 0.005   | 0.001      | − 9.537    | 0.000    |
| *Population size (contrast: 100)* | |          |            |          |
| 101                      | − 0.008   | 0.001      | − 15.579   | 0.000    |

with polarization across boundary values. However, polarization decreases when $\mu_{SPIRO}$ is raised from 0.61 to 0.73 and 0.85. We interpret this as a consequence of the dominant system dynamics transitioning from polarized state to fractured state for the highest $\mu_{SPIRO}$ values (see Sect. 4).

**Fig. 2** ESBG-$\mu_\varepsilon$ relationship for model VBVI for different values of $\mu_{SPIRO}$. Panels represent different values of $\sigma_\varepsilon$



**Fig. 3** ESBG-$\mu_{SPIRO}$ relationship for model VBVI for different values of $\mu_\varepsilon$. Panels represent different values of $\sigma_\varepsilon$

## 4 Discussion and Future Work

In this work we implemented a novel algorithm for dynamic detection of identity groups based on their opinions. In recognition of the common observation that people differ in their judgements on how many partisan groups there are in a society, and which individual belongs to which group, we parameterized our implementation of identity with the variable we call SPIRO.

SPIRO determines how closely a pair of agents must be to be considered for identity group detection. Through visual inspection of the course of the models' runs, it appears that higher SPIRO values (0.73 and 0.85) causes the opinion space to be split into more identity groups. The effects of these parameters will be explored in detail elsewhere (in preparation). In our last model we allow SPIRO to vary across agents to account for people perceiving different sets of identity groups around them. This makes our model more realistic, while being computationally efficient due to our method of partitioning.

Through our analysis of the behaviors of our models, we are able to determine which experimental variables in our different models are the most relevant for polarization. We find that models with identity exhibit a higher average polarization across their different experimental conditions than models without identity. We also find that introducing heterogeneity in both Boundary and SPIRO in our model lowers polarization overall. This is admittedly a simplistic way of analyzing the effects of identity and heterogeneity. We will dive deeper into the role of these model features in a future article.

We also observe that the influence of identity on polarization depends on the SPIRO value of the agents. For moderate values of mean SPIRO, polarization monotonically increases with SPIRO. However, the highest two SPIRO values we have considered here show a deviation from this trend and show reduced polarization. Since the ESBG algorithm privileges bi-polarization over fractured states with multiple tight clusters, this can be explained by a fracturing of the agent population into several opinion camps. This is another aspect of our analysis that we will discuss in more detail in a future work.

Going forward, we will also be looking at the effect of heterogeneity of boundary and SPIRO on the behavior of the system. Previous studies have looked at the influence of boundary heterogeneity on consensus [20] and the number and size of opinion clusters [6, 21]. Consistent with these studies we find that heterogeneous $\varepsilon$ causes the system to tend towards less polarized states, possibly towards consensus. This is likely due to the possibility that agents with above-average $\varepsilon$ act as bridging agents due to their openness to a wider range of opinions, while the agents with below-average $\varepsilon$ might not have much of an influence on the system dynamics. We purport a similar mechanism might be at play in the case of the heterogeneous SPIRO model— variance in group classification might lead to less clearly defined identity bubbles, which would allow some agents to act as bridges between clusters that emerge due to identity effects.

# References

1. Hegselmann, R., Krause, U.: Opinion dynamics and bounded confidence models, analysis and simulation. J. Artif. Soc. Soc. Simul. **5**, 1–2 (2002)
2. Lanchier, N., Li, H.-L.: Consensus in the Hegselmann-Krause model. J. Stat. Phys. (2022). https://doi.org/10.1007/s10955-022-02920-8
3. Kurz, S.: How long does it take to consensus in the Hegselmann-Krause model? PAMM 803–804 (2014). https://doi.org/10.1002/pamm.201410382
4. Matakos, A., Terzi, E., Tsaparas, P.: Measuring and moderating opinion polarization in social networks. Data Mining Knowl. Discov. 1480–1505 (2017). https://doi.org/10.1007/s10618-017-0527-9
5. Pineda, M., Toral, R., Hernández-García, E.: The noisy Hegselmann-Krause model for opinion dynamics. Eur. Phys. J. B (2013). https://doi.org/10.1140/epjb/e2013-40777-7
6. Fu, G., Zhang, W.: Opinion dynamics of modified Hegselmann-Krause model with group-based bounded confidence. IFAC Proc. Volumes 9870–9874 (2014). https://doi.org/10.3182/20140824-6-za-1003.02770
7. Nedic, A., Touri, B.: Multi-dimensional Hegselmann-Krause dynamics. In: 2012 IEEE 51st IEEE Conference on Decision and Control (CDC) (2012). https://doi.org/10.1109/cdc.2012.6426417
8. Parasnis, R., Franceschetti, M., Touri, B.: Hegselmann-Krause dynamics with limited connectivity. In: 2018 IEEE Conference on Decision and Control (CDC) (2018). https://doi.org/10.1109/cdc.2018.8618877
9. Yixuan, D., Cheng, T., Shing, W.W.: Discrete-time Hegselmann-Krause model for a leader-follower social network. In: 2018 37th Chinese Control Conference (CCC) (2018). https://doi.org/10.23919/chicc.2018.8482680
10. Hegselmann, R., Krause, U., et al.: Opinion dynamics under the influence of radical groups, charismatic leaders, and other constant signals: a simple unifying model. Netw. Heterog. Media 477–509 (2015). https://doi.org/10.3934/nhm.2015.10.477
11. Tajfel, H., Turner, J.C., Austin, W.G., Worchel, S.: An integrative theory of intergroup conflict. Organ. Identity Reader **56**, 9780203505984–9780203505916 (1979)
12. Baldassarri, D., Gelman, A.: Partisans without constraint: political polarization and trends in American public opinion. AJS **114**, 408–446 (2008)
13. Iyengar, S., Lelkes, Y., Levendusky, M., Malhotra, N., Westwood, S.J.: The origins and consequences of affective polarization in the United States. Annu Rev Polit Sci. **22**, 129–146 (2019)
14. Wojcieszak, M., Kelly Garrett, R.: Social identity, selective exposure, and affective polarization: how priming national identity shapes attitudes toward immigrants via news selection. Human Commun. Res. 247–273 (2018). https://doi.org/10.1093/hcr/hqx010
15. Slater, M.D.: Reinforcing spirals: the mutual influence of media selectivity and media effects and their impact on individual behavior and social identity. Commun. Theory (2007). https://academic.oup.com/ct/article-abstract/17/3/281/4098761
16. Slater, M.D.: Reinforcing spirals model: conceptualizing the relationship between media content exposure and the development and maintenance of attitudes. Media Psychol. **18**, 370–395 (2015)

17. Blondel, V.D., Guillaume, J.-L., Lambiotte, R., Lefebvre, E.: Fast unfolding of communities in large networks. J. Stat. Mech. **2008**, P10008 (2008)
18. NetLogo code, data, and analysis scripts, https://github.com/frantisek901/Spirals/tree/master/IdentityPaper
19. Tang, T., Ghorbani, A., Squazzoni, F., Chorus, C.G.: Together alone: a group-based polarization measurement. Qual. Quant. **56**, 3587–3619 (2022)
20. Chen, G., Su, W., Ding, S., Hong, Y.: Heterogeneous Hegselmann-Krause dynamics with environment and communication noise. IEEE Trans. Automat. Contr. **65**, 3409–3424 (2020)
21. Han, W., Huang, C., Yang, J.: Opinion clusters in a modified Hegselmann–Krause model with heterogeneous bounded confidences and stubbornness. Phys. A Stat. Mech. Appl. 121791 (2019). https://doi.org/10.1016/j.physa.2019.121791

# Influence and Similarity in Social Networks: A Study of the Opinion Dynamics Among Teenagers Through an Agent-Based Model

**Dario Germani** [ID]**, Barbara Sonzogni** [ID]**, and Federico Cecconi** [ID]

**Abstract** The main goal of this research is to study the dynamics of the opinion among teenagers by reconstructing the processes of influence that take place during their interactions, raising further questions about the ways and reasons why individuals get in touch with others. The integration between Agent-Based Modelling and sociometry allowed the conceptualization of the phenomenon as a diffusion study, considered as the outcome of the imitative process triggered by any compliance motives especially in view of the sociological tradition. In particular, the concepts of social influence and homophily can be traced back to a dual mechanism able of explaining it: (1) the behavior of peers occupying a relevant position within relational groups (school classes); (2) the interaction favored by specific elements linked to the similarity between individuals. The empirical results obtained from a web survey has been compared with the ones from the simulation model in order to reproduce the above social phenomenon and to confirm the theoretical assumptions behind the model itself.

D. Germani (✉)
Department of Political and Social Sciences, University of Cagliari, Via Fra Ignazio 78, 09123 Cagliari, Italy
e-mail: dario.germani@unica.it

B. Sonzogni
Department of Communication and Social Research (CORIS), Sapienza University of Rome, Via Salaria 113, 00198 Rome, Italy
e-mail: barbara.sonzogni@uniroma1.it

F. Cecconi
LABSS (Laboratory of Agent Based Social Simulation), Institute of Cognitive Sciences and Technologies (CNR), via Palestro 32, 00185 Rome, Italy
e-mail: federico.cecconi@istc.cnr.it

# 1 Introduction

## *1.1 Theoretical Framework*

Music at any ages is an essential component that marks the important steps of people's lives: it drives the beginning or the end of a love affair or a friendship, it helps people to process and overcome a disappointment, or it revives good times. The sociology of music has received a significant interest in recent decades. Early work in this area had already raised important questions, such as how individuals listen to the music [1, 2], how music figures in community life [3] and how musical preferences vary within a population by documenting the patterns and preferences in different audiences [4]. Taste and competence deserve instead a separate discourse in terms of cultural capital: if Bourdieu [5] argued that artistic preferences and competences are class-based, today traditional status distinctions are becoming more and more faded as far as the emergence of a *cultural omnivore* style of consumption [6].

Although music reception may seem a private and isolated activity, many studies point out that it is also a group activity since people with similar choice options may gravitate towards each other. Such studies tend to gloss over the music content, emphasising the lifestyle elements that distinguish music-based groupings such as fashion and public behavior [7]. According to research conducted by North et al. [8], teenagers spend between two and three hours a day listening to the music. It is prevalent enough in their teens to intrigue those authors who have attempted to justify this practice by defining two main functions, one like solo and the other one like collective. Considering the last one, music can be often at the root of the formation of groups or the reason for belonging to one of them: people tend to get closer and create relationships with their peers and, in addition to similarity of trend, likenesses may occur in musical choices.

This leads to the following questions: (i) what kind of role do opinions on musical preferences play in social interactions? (ii) How do actors select each other building dynamically their relational structure over time?

Even if social encounters are reckoned us spontaneous events, scientific research affirms that *sympathy*, *closeness* and *similarity* are key points in the formation of social relationships [9–14]. This set of elements we usually take for granted are not visible on the surface: they lie below the observation plane influencing the interaction process among individuals.

Within this framework, in the first place studying musical opinions is expressing what they contribute to social relations.

## *1.2  Opinion Dynamics*

Opinion formation is influenced by the combination of several elements: self-reflection, external sources of information and real-world experiences that supports the individual's reasoning process. Moreover, social interactions within the communication system play a critical role in their formation.

Agent-Based Modelling for the study of this topic has become an independent area in social sciences with a strong interdisciplinary attention. Existing representations of opinion dynamics can be divided on two different features[1]: (a) the representation of the opinions through binary, discrete or continuous variables; (b) the local rules of interactions between agents that reflect the basic theories of the models. The most relevant are:

- the *voter model* [15], which investigates the trend of public preferences, such as the voting choices between two candidates; its essential version (*peer-to-peer*) involves a set of agents, in a defined space, whose opinions may change in relation to their neighbours, giving effect to a global process connecting the entire population; another type is the so-called *majority model* [16], where agents try to follow the mainstream;
- the *Sznajd model* [17], which is based on the assumption that it is easier to be convinced by two or more agents sharing the same opinion than to be convinced by only one;
- the *dissemination of culture* [18], based on the assumptions that agents are more likely to interact with the ones who share many of their features (1); these exchanges tend to increase the number of features they have in common, thus increasing the likelihood of interacting again (2); according to the author, *culture* represents the set of individual traits susceptible to the social influence;
- the *bounded confidence* [19, 20], according to which each agent has an opinion that can change when it becomes conscious of the opinions of its neighbours within a relational network. In order to be mutually influenced, two interacting agents must have opinions close enough: if the difference between of the opinions is significant, the communication process is impossible and there will be no change in their respective views; the result of the interaction is a sort of compromise towards the other's opinion unless their dissension is below a given threshold.

The purpose of these models is to combine theoretical assumptions with hypotheses described in algorithmic terms through computer simulations in order to explain and replicate the formation process of the opinions.

---

[1] The following review cannot do justice to all the contributions in the literature, but we believe that our classification can provide as a general guideline for the development and communication of the opinion formation models.

**Table 1** Conceptualised features grid

| Students features | Network | Musical opinions |
|---|---|---|
| Gender | Nodes | Rock |
| School performance | Edges | Classic |
| Parents education | | Jazz |
| Music competence | | Dance |
| Music enjoyment | | Pop |
| | | Rap/Hip hop |
| | | Trap |
| | | Reggae |
| | | Indie |

## 2   Method

### 2.1   Web Survey

The collection of the information was carried out through the construction and administration of a self-compiled online questionnaire divided into 24 questions for a total of 85 respondents, divided into seven classes of two different schools. Having interviewed limited groups, such as the students of the school classes in the present case, it has been achieved the use of a proposal of integration between web survey and sociometry, thus deepening the structure and the intensity of the relationships among classmates.

As originally intended [21] and according to certain principle of choice, the sociometric test concerns the preferences of each member of a group for the other members of the same group. In this work, we decided to focus on 5 possible choices notifying the possibility of expressing a lower number, also weighting these combinations according to a preference ranking (*weight*): the first option corresponds to a score of 5 points, the second one to a score of 4 points, and so on. As a result of this scheme, it has been possible reconstructing the school classes social maps and identifying the presence of central nodes.

A variety of aspects concerning modes of music consumption/performance, opinions on different musical genres and socio-demographic individual features have been investigated, too (Table 1).

### 2.2   Web Survey

The model integrates some features of the *bounded confidence* [19, 20] and the *dissemination of culture* [18] model from real data assuming a hybrid interaction

dynamic: agents influence each other with a strength based on a *convergence* parameter ($\mu$), which tells how strongly the agent $x$ gets the opinion of the agent $y$, but only if the difference between the features of the agents is less than a *tolerance* parameter ($\theta$).[2] Considering a population of $N$ agents, the opinion space is $[a, b] \in R$. If $x$ and $y$ agents are randomly selected and they meet at a $t$ time with opinions $[a, b] \in R$, the interaction rule is as follows:

$$(\eta_t(x), \eta_t(y)) = ((a + \mu(b - a), b + \mu(a - b))$$
$$if \ [a - b] \leq \theta(a, b) \tag{1}$$

where $\eta_t(x)$ is the opinion of $x$ agent at a $t$ time.

On one hand, these two inner workings suppose the exchange of views takes place according to the *similarity* among the agents. On the other hand, if the theory of *social influence* [22–25] is developed within a more relational paradigm, has been possible to hypothesise the presence of agents within a relational network with different degrees of expertise even in the music field (defined as *leadership opinion*) based on certain features.

The distance among the features of the agents (*similarity*) is a real number between zero and one (2): where the more the agents are similar the more the similarity value is close to zero; the more they are different the more the similarity value is close to one.

$$D(x, y) = \frac{[(x_1 - y_1) + (x_2 - y_2) + \cdots (x_n - y_n)]}{n} \tag{2}$$

The simulation is set up in the following way: we have an agent-set (students) reproducing the same features of those obtained from the empirical analysis (i); a network topology reproducing the same structure obtained from the empirical analysis (ii); a spread of music opinions[3] through a random algorithm at the first step (iii). Based on this, we wonder whether it will be possible to reproduce a similar or a different opinion dynamic at the final step as empirically observed.[4]

---

[2] The model uses the *tolerance* and *convergence* parameter—both with a value between zero and one—where the former is considered as a similarity-based confidence threshold describing an agent's resistance to alternative points of view. If the difference between the features of the two agents is below it, the gap can be reduced by reaching a kind of a compromise of one or the other, otherwise they keep their current opinions after the interaction. The latter measures instead the influence capacity of the model, as a multiplier stating the relative agreement between the involved agents.

[3] The answers on each genre have been recoded into a discrete variable from 1 to 4 (1 = *strongly like it*; 2 = *like it*; 3 = *don't like it*; 4 = *don't listen it*).

[4] The analyzed interaction dynamics will focus on the students' favourite music genre (Rap/Hip Hop).

**Table 2** Average value of the distance of opinion as a function of the similarity of interactions: the table shows the results at the aggregate level

| | Similarity | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | 0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 |
| Mean opinion distance value (web survey) | 0.15 | 0.66 | 0.42 | 0.75 | 1.15 | 1.15 | 1.13 | 1.7 | 1 |

**Table 3** Average value of the distance of opinion as a function of the intensity of ties: the table shows the results at the aggregate level

| | Weight | | | | |
|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 |
| Mean opinion distance value (web survey) | 1.1 | 1.12 | 0.9 | 0.83 | 0.83 |

**Table 4** Illustrative overview of the opinion dynamics at the micro level

| Talker | Receiver | Similarity | Weight | Opinion distance |
|---|---|---|---|---|
| A | B | 0.3 | 5 | 1 |
| A | C | 0.5 | 4 | 2 |
| A | D | 0.1 | 3 | 0 |
| B | A | 0.4 | 1 | 1 |
| … | … | … | … | … |
| X | Y | 0.2 | 1 | 0 |

# 3   Results

## 3.1   Empirical Analysis

The idea that people tend to interact more with their own kind has been empirically validated also in the course of this investigation: in fact, the average value of the opinion distance increases as the similarity of interactions decreases and as the intensity of ties of interactions decreases: we can observe an increasing trend if we consider the similarity among the students (Table 2) and a decreasing trend if we consider the intensity of the relationships, instead (Table 3). Basically, students with the same opinion on a certain musical genre are also those characterised by high similarity and strong ties.

The tables as shown above are the aggregate-level[5] representation of the following pattern at the micro level (Table 4).

---

[5] Data analysis has been performed on a total of 221 interactions on the *sympathy* network topology, the sociometric dimension through which we asked to the respondents the following question: "among your classmates, who do you like the most?".

**Table 5** Average value of the distance of opinion as a function of the similarity of interactions

|  | Similarity | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
|  | 0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 |
| Mean opinion distance value ($\mu$ & $\theta$ = 0.1) | 0.92 | 1.09 | 1.04 | 1.29 | 1.31 | 1.37 | 1.15 | 1.66 | 1.33 |
| Mean opinion distance value ($\mu$ & $\theta$ = 0.2) | 1.07 | 1.09 | 0.95 | 1.29 | 1.24 | 1.33 | 1.18 | 1.22 | 1.33 |
| Mean opinion distance value ($\mu$ & $\theta$ = 0.3) | 1.07 | 1.09 | 0.91 | 1.26 | 1.24 | 1.21 | 1.18 | 1.44 | 1.33 |
| Mean opinion distance value ($\mu$ & $\theta$ = 0.4) | 1.07 | 0.9 | 1 | 1.32 | 1.34 | 1.21 | 1.21 | 1.55 | 1.33 |
| Mean opinion distance value ($\mu$ & $\theta$ = 0.5) | 1 | 1.09 | 0.91 | 1.23 | 1.24 | 1.13 | 0.9 | 1.66 | 1.33 |
| Mean opinion distance value ($\mu$ & $\theta$ = 0.6) | 1.07 | 1.09 | 0.95 | 1.26 | 1.24 | 1.13 | 1.21 | 1.44 | 1.33 |
| Mean opinion distance value ($\mu$ & $\theta$ = 0.7) | 0.46 | 0.27 | 0.79 | 0.76 | 0.96 | 0.94 | 0.93 | 1.88 | 1.66 |
| Mean opinion distance value ($\mu$ & $\theta$ = 0.8) | 0.38 | 0.54 | 0.79 | 0.76 | 1.17 | 1.03 | 0.93 | 1.88 | 1.66 |
| Mean opinion distance value ($\mu$ & $\theta$ = 0.9) | 0.15 | 0.45 | 0.37 | 0.32 | 0.65 | 0.68 | 0.78 | 1.88 | 1.66 |
| Mean opinion distance value ($\mu$ & $\theta$ = 1) | 0.23 | 0.81 | 0.25 | 0.29 | 0.55 | 0.54 | 0.51 | 1.11 | 1 |

## 3.2 Simulative Results

The NetLogo software has been set up in order to run the model many times with the chance to change its settings and recording the results of each run.[6] This process allows us to explore different configurations and different behaviors in the action system: each time unit (*step*) corresponds to an execution of the influence process where each agent connects with its partners.

To understand what is the best match among parameters capable of reproducing an opinion dynamic more or less similar to the one empirically observed, we implemented the model executing some initial experiments in view of numerous scenarios: due to this, it has been possible to examine the robustness of the model outcomes[7] in relation to the changes of the parameter values.

The simulation allowed us to generate data that were perfectly comparable with the empirical ones: varying the parameters of the theorised mechanisms by a value of 0.1 each time, ten different scenarios have been generated (Tables 5 and 6).

---

[6] Model, code and results are available at the following link: https://github.com/DarioGermani/Music-opinion-dynamic.

[7] Sensitivity analysis aims to ascertain interaction effects by sampling the model output over a wide range of parameter values [26].

**Table 6** Average value of the distance of opinion as a function of the intensity of ties

|  | Weight | | | | |
|---|---|---|---|---|---|
|  | 1 | 2 | 3 | 4 | 5 |
| Mean opinion distance value ($\mu$ & $\theta = 0.1$) | 1.25 | 1.36 | 1.2 | 1.07 | 1.35 |
| Mean opinion distance value ($\mu$ & $\theta = 0.2$) | 1.07 | 1.13 | 1.3 | 1.19 | 1.28 |
| Mean opinion distance value ($\mu$ & $\theta = 0.3$) | 1.18 | 1.11 | 1.37 | 0.98 | 1.28 |
| Mean opinion distance value ($\mu$ & $\theta = 0.4$) | 1.25 | 1.19 | 1.35 | 1.03 | 1.28 |
| Mean opinion distance value ($\mu$ & $\theta = 0.5$) | 0.96 | 1.11 | 1.37 | 1 | 1.13 |
| Mean opinion distance value ($\mu$ & $\theta = 0.6$) | 1.22 | 1.02 | 1.32 | 1.01 | 1.28 |
| Mean opinion distance value ($\mu$ & $\theta = 0.7$) | 1.03 | 1.05 | 0.92 | 0.88 | 0.66 |
| Mean opinion distance value ($\mu$ & $\theta = 0.8$) | 1.14 | 1.11 | 0.92 | 0.96 | 0.73 |
| Mean opinion distance value ($\mu$ & $\theta = 0.9$) | 0.88 | 0.88 | 0.57 | 0.49 | 0.47 |
| Mean opinion distance value ($\mu$ & $\theta = 1$) | 1.03 | 0.63 | 0.47 | 0.31 | 0.3 |

## 4 Conclusion

In the beginning, we identified a common approach to replicate the dynamic of the opinions among teenagers through ABMs in literature. On one hand, this research represents the attempt to test the predictions of theoretical assumptions implemented into a simulative model: it was not just a matter of observing the phenomenon but reproducing it by doing experiments; on the other hand, it contributes to filling a gap to the relation between music, opinion formation process and social interactions.

Following the lead of *bounded confidence* and *dissemination of culture*, we hypothesized that the dynamic through which people interact and influence each other in this sphere (and we may presume in others as well) has to do with similarity. The model outcomes show that many scenarios described accurately the same pattern empirically observed, in particular when we increase the value of the parameters. The most important result is the average value of the opinion distance among the students that increases as the similarity decreases and as the intensity of ties decreases.

Therefore, the algorithm we suggested poses how individuals interact and convey opinions on musical preferences: based on this, Agent-Based Modelling has proved to be a powerful tool for testing hypotheses on social mechanisms.

# References

1. Riesman, D.: Listening to popular music. Am. Q. **2**(4), 359–371 (1950)
2. Hatch, D.J., Watson, D.R.: Hearing the blues: an essay in the sociology of music. Acta Sociologica **17**(2), 162–177 (1974)
3. Coleman, J.S.: The Adolescent Society. Free Press of Glencoe (1961)
4. Lazarsfeld, P.F., Stanton, F.: Radio Research. Sloan and Pearce, New York, Duell (1941)
5. Bourdieu, P.: Distinction: A Sociological Critique of the Judgement of Taste. Cambridge University Press, Cambridge, UK (1984)
6. Peterson, R.A., Kern, R.M.: Changing highbrow taste: from snob to univore. Am. Sociol. Rev. **61**(5), 900–907 (1996)
7. Weinstein, D.: The sociology of rock: an undisciplined discipline. Theory Cult. Soc. **8**(4), 97–109 (1991)
8. North, A.C., Hargreaves, D.J., O'Neill, S.A.: The importance of music to adolescents. Br. J. Educ. Psychol. **70**(2), 255–272 (2000)
9. Festinger, L., Schachter, S., Back, K.: Social Pressures in Informal Groups: A Study of Human Factors in Housing. Stanford University Press, Stanford (1950)
10. Lazarsfeld, P.F., Merton, R.K.: Friendship as a social process: a substantive and methodological analysis. Freedom Control Mod. Soc. **18**(1), 18–66 (1954)
11. Byrne, D., Clore, J.L., Worchel, P.: Effect of economic similarity-dissimilarity on interpersonal attraction. J. Pers. Soc. Psychol. **4**(2), 220–224 (1966)
12. Osbeck, L.M., Moghaddam, F.M., Perreault, S.: Similarity and attraction among majority and minority groups in a multicultural context. Int. J. Intercult. Relat. **21**(1), 113–123 (1997)
13. Drigotas, S.M.: Similarity revisited: a comparison of similarity—attraction versus dissimilarity—repulsion. Br. J. Soc. Psychol. **32**(4), 365–377 (1993)
14. Schaefer, D.R., Simpkins, S.D., Vest, A.E., Price, C.D.: The contribution of extracurricular activities to adolescent friendships: new insights through social network analysis. Dev. Psychol. **47**(4), 1141–1152 (2011)
15. Clifford, P., Sudbury, A.: A model for spatial conflict. Biometrika **60**(3), 581–588 (1973)
16. Galam, S.: Real space renormalization group and totalitarian paradox of majority rule voting. Phys. A **285**(1–2), 66–76 (2000)
17. Sznajd, W.K., Sznajd, J.: Opinion evolution in closed community. Int. J. Modern Phys. C **11**(06), 1157–1165 (2000)
18. Axerlord, R.: The dissemination of culture: a model with local convergence and global polarization. J. Confl. Resol. **41**(2), 203–226 (1997)
19. Deffuant, G., Neau, D., Amblard, F., Weisbuch, G.: Mixing beliefs among interacting agents. Adv. Complex Syst. **3**, 87–98 (2000)
20. Deffuant, G., Amblard, F., Weisbuch, G., Faure, T.: How can extremism prevail? A study based on the relative agreement interaction model. J. Artif. Soc. Social Simul. 5(4) (2002). https://www.jasss.org/5/4/1.html
21. Marineau, R.F.: The birth and development of sociometry: the work and legacy of Jacob Moreno (1889–1974). Social Psychol. Q. **70**(4), 322–325 (2007)
22. Katz, E., Lazarsfeld, P.F.: Personal Influence: The Part Played by People in the Flow of Mass Communications. Free Press, New York (1955)
23. Akers, R.L., Krohn, M.D., Lanza-Kaduce, L., Radosevich, M.: Social learning and deviant behavior: a specific test of a general theory. Am. Sociol. Rev. **44**(4), 636–655 (1979)
24. Asch, S.E.: Studies of independence and conformity: I. a minority of one against a unanimous majority. Psychol. Monogr. General Appl. **70**(9), 1–70 (1956)
25. Myers, D.G., Lamm, H.: The group polarization phenomenon. Psychol. Bull. **83**(4), 602–627 (1976)
26. Saltelli, A., Tarantola, S., Campolongo, F., Ratto, M.: Sensitivity Analysis in Practice: A Guide to Assessing Scientific Models. England, Chichester (2004)

# Networked Models of Social Influence: Explaining Left-Right Political Landscapes in Europe Through Opinion Dynamics and Network Structure

**Daniel Reisinger, Michael Vogrin, Guilherme Wood, Thomas Schmickl, and Georg Jäger**

**Abstract** Traditional models of social influence typically use assimilative or repulsive influence to study how consensus or polarization emerge. Given simple network structures, such as fully connected graphs, traditional models often fail to account for the multi-modal opinion distributions found in empirical data. In this study, we focus on more realistic social network structures in terms of clustering coefficient and average shortest path length and construct a model that allows both assimilative and repulsive influence to drive opinion changes in individuals. We find that non-trivial patterns emerge when the forces of assimilative and repulsive influence are kept at a specific ratio and the network structure is highly clustered. Comparisons with empirical left-right political opinion landscapes show that our model produces realistic results that share the multi-model characteristics as observed in data collected by the European Social Survey program.

**Keywords** Social influence · Opinion dynamics · Polarization · Political landscapes · Networks · Micro-macro link

## 1 Introduction

Social influence can be defined as change in an individual's opinion, attitude, feeling, thought, or behavior in response to interactions with other individuals and encompasses the mechanisms that drive this change [13]. Resulting change of opinions

D. Reisinger (✉) · M. Vogrin · G. Wood · T. Schmickl · G. Jäger
University of Graz, Graz, Austria
e-mail: daniel.reisinger@uni-graz.at

D. Reisinger · G. Jäger
Institute of Environmental Systems Sciences, Graz, Austria

M. Vogrin · T. Schmickl
Institute of Biology, Graz, Austria

G. Wood
Institute of Psychology, Graz, Austria

can thus be conceptualized as a function of that individual's opinion, the opinions of its surrounding, and the kind of social influence, which ultimately determines the direction of the opinion change. Literature identifies many different types of social influence including conformity, obedience, persuasion, and differentiation [11]. Conformity, obedience, and persuasion clearly point in the direction of reducing opinion differences, whereas differentiation points towards increasing opinion differences. The idea of conformity can be used to explain social stability through convergence of opinions [16]. Analogously, the idea of differentiation can be used to explain the phenomenon of polarization [17]. To understand the dynamics that arise from the many different types of social influence, numerous models have been proposed to capture the conditions under which phenomena such as opinion convergence, polarization, or clustering emerge.

In a recent review of models of social influence, it is argued that much of the literature on social influence and opinion dynamics revolves around three classes of models: (i) models of assimilative influence, (ii) models with similarity biased influence, and (iii) models with repulsive influence [5]. The key assumptions behind these models as well as their core results are briefly summarized here.

(i) Models of assimilative social influence assume that individuals with different opinions move towards reducing these opinion differences [5]. This assumption is supported by classical psychological findings such as Asch's conformity experiment [2], the effects of peer influence [1], theoretical models such as balance theory [8], and frameworks that focus on persuasion [12]. Assimilative social influence typically leads to consensus among individuals in the long run [5].

(ii) Models of similarity biased influence assume that social influence is dependent on the degree of difference in opinions of individuals such that only sufficiently similar individuals can influence each other towards reducing opinion differences [5]. The key assumption of similarity biased influence can be supported, for example, by the "Backfire Effect", in which individuals are not persuaded by facts that contradict their beliefs, but instead fortify their positions [15]. Models of similarity biased influence typically lead to consensus if the similarity bias is low, but may also show patterns of opinion clustering if the similarity bias is high [5, 6].

(iii) Models of repulsive influence draw on the assumption that if individuals are too dissimilar they may influence each other towards increasing their opinion differences [5]. Salzarulo justifies repulsive influence by referring to self-categorization theory [14]. Others point to specific phenomena that are tied to repulsive influences such as xenophobia or differentiation [10]. Generally, models of repulsive influence predict that individuals form opposing opinion groups resulting in patterns of bi-polarization.

The question remains how useful these models are in replicating and explaining empirical data such as shown in Fig. 1 [9]. In this regard, it has been argued that (a) the central peaks in such data suggest the presence of assimilative social influence, (b) the non-central opinion clusters can be taken as evidence of similarity-biased

**Fig. 1** Left-right political landscapes. Respondents from Europe (aggregated) and selected European countries placed themselves politically on a left-right spectrum in 2018. Data are obtained from the European Social Survey [4]

social influence, and (c) the extremal peaks on the far left and far right in the opinion spectrum indicate repulsive influences in the social system [5]. Thus, the models of social influence described above appear to be capable of capturing the core characteristics of empirical left-right political opinion landscapes. However, they do so only in isolation, meaning one characteristic at a time. The model classes fail to capture the pattern in its entirety, i.e. they fail to reproduce the multi-modal opinion distributions shown in Fig. 1. In this study, we continue the works by Flache et al. and study how one might achieve a better model fit without additional model assumptions on the social influence mechanisms. Instead, we want to focus on how the empirical left-right opinion landscape can be explained through a balancing of assimilative and repulsive forces in combination with more realistic network structures.

## 2 Methodology

### 2.1 Social Influence Mechanisms

Social influence is modelled following traditional opinion dynamics equations. Individuals are represented by nodes in a network where every node has an opinion $o_i$ in the closed interval [0, 1]. The interval boundaries represent an extreme left and an extreme right political opinion and a value of 0.5 represents neutrality or indifference. The structural relationship between individuals is described by a network where individuals may only interact with one another if they are directly connected, i.e. immediate neighbors. The equations describing the social influence mechanisms

are modelled following [5]. In this regard, we distinguish between an assimilative mechanism (with similarity bias) and a repulsive mechanism. The general updating equation for either mechanism follows

$$o_{i,t} = o_{i,t-1} + \mu \sum_j f_w \cdot (o_{j,t-1} - o_{i,t-1}) \tag{1}$$

where $o_i$ represents the opinion of node $i$ and $o_j$ the opinion of one of node $i$'s neighbors. The parameter $\mu$ represents the rate of convergence of a particular social influence mechanism and the function $f_w$ describes the working principle of the social influence mechanisms itself. For assimilative influence, the function $f_w^a$ is defined as

$$f_w^a = \begin{cases} 1, & \text{if } |o_j - o_i| < \epsilon. \\ 0, & \text{otherwise.} \end{cases} \tag{2}$$

where $\epsilon$ describes a confidence level that determines how close two opinions have to be for assimilative influence to take effect. For repulsive influence, the function $f_w^r$ is defined as

$$f_w^r = 1 - \omega \cdot |o_j - o_i| \tag{3}$$

where $\omega$ is used to set the critical level of opinion difference for the repulsive influence to take effect. Note that for the above equation the opinions are not naturally bound on the opinion interval [0, 1] and require occasional truncation [5]. Figures 2 and 3 show example opinion dynamics for assimilative influence with varying parameter values of $\epsilon$ and for repulsive influence with varying parameter values of $\omega$ on a complete graph. Opinions are initialized randomly and uniformly over the interval [0, 1]. Individuals interact based on the social influence mechanisms defined above where every line represents the dynamic of a single individual's opinion over time. Depending on the value of $\epsilon$, assimilative influence can create opinion clustering as well as full consensus (see Fig. 2). Depending on the value of $\omega$, repulsive influence can create full consensus as well as bi-polarization (see Fig. 3).

We then construct a model to account for both social influences including assimilative and repulsive dynamics where the updated opinion of an individual is simply the combined change in opinion of the assimilative and the repulsive mechanisms as



**Fig. 2** Assimilative mechanism (with similarity bias). Shows opinion dynamics for different values of $\epsilon$ on a complete graph with 100 nodes

**Fig. 3** Repulsive mechanism. Shows opinion dynamics for different values of $\omega$ on a complete graph with 100 nodes

outlined above. The updated opinion of an individual in the combined social influence setting is calculated following

$$o_{i,t} = o_{i,t-1} + \Delta^a + \Delta^r \tag{4}$$

$$\Delta^a = \mu^a \sum_j f_w^a \cdot (o_{j,t-1} - o_{i,t-1}) \tag{5}$$

$$\Delta^r = \mu^r \sum_j f_w^r \cdot (o_{j,t-1} - o_{i,t-1}) \tag{6}$$

In contrast to the assimilative mechanism (Eq. 2) which may only produce opinion convergence, and the repulsive mechanism (Eq. 3) which, on its own, is already able to reproduce both opinion convergence and opinion polarization, the combined model should be interpreted as follows: Every individual in the system experiences assimilative influence due to the assimilative mechanism (Eq. 2) which gets, depending on local network structure, further amplified or reduced due to the addition of the repulsive mechanism (Eq. 3).

## 2.2 Watts-Strogatz Networks

Social network structure plays an important part in the study of opinion dynamics as the properties of a specific network may accelerate or decelerate the opinion change induced by the mechanisms of social influence. Selecting a network structure with similar properties as observed in real social networks is, thus, crucial to the assessment of different social influence mechanisms. For this purpose, we select a graph generating method known as the Watts-Strogatz model [18] which allows us to vary two network properties especially relevant to social networks: (i) the average clustering coefficient $C$ and (ii) the average shortest path length $L$ [3]. The properties may be varied by changing the model's rewiring probability $p$, which determines the probability of selecting and rewiring an edge in an initially highly clustered graph. This allows us to generate vastly different networks in terms of average clustering

**Fig. 4** Properties of
Watts-Strogatz graphs.
Shows the average clustering
coefficient $C$ and the average
shortest path length $L$ for
different rewiring
probabilities $p$ in the
Watts-Strogatz graph with
average degree 10 and 200
nodes



coefficient and average shortest path length. Note that real social networks typically exhibit a high average clustering coefficient and a low average shortest path length [3]. In Fig. 4, this would correspond to networks generated with rewiring probability $p \approx 0.01$.

## 2.3 Simulation Setup

To capture the opinion patterns produced by the combined social influence model on different network structures we devise the following simulations setup: We generate Watts-Strogatz networks with 200 nodes where every node is randomly and uniformly assigned an opinion value between 0 and 1. The opinion of a node $o_i$ is updated following Eq. 4 where a node's neighboring opinions $o_j$ are given by the structural relationship of the network. All opinions are updated in random sequential order and the updating procedure is repeated for a simulation period of 200 time steps. The stopping condition of 200 time steps is oriented on the drop of average opinion differences between timesteps—$1/N \sum_{i=1}^{N} |o_{i,t-1} - o_{i,t}|$—which shows a significant slowing down of the system after an initial 100 timesteps, thus, signifying a relatively stable phase. At 200 time steps the system's average opinion difference between two timesteps is already well below $10^{-3}$. At the end of the simulation we take the opinion distribution over all nodes in the form of a histogram with eleven equal sized bins, corresponding to the bin size of the empirical data shown in Fig. 1.

To identify (i) the space of non-trivial patterns and (ii) the model fit with empirical data, we perform a parameter sweep over the rates of convergence, $\mu^a$ and $\mu^r$, driving the change of the assimilative and repulsive mechanism. Both parameters are varied iteratively in 30 steps over the closed interval $[0, 0.1]$ while the confidence level of assimilative influence is fixed at $\epsilon = 1$ and the critical level of repulsive influence is fixed at $\omega = 2$. For assimilative influence, $\epsilon = 1$ corresponds to full

opinion convergence of structurally linked nodes (see Fig. 2) excluding any form of similarity bias. For repulsive influence, $\omega = 2$ corresponds to a critical level for opinion differences of 0.5, meaning that any opinion difference between two nodes greater than 0.5 leads to a repulsive influence between those nodes and any opinion difference smaller than 0.5 to an assimilative influence between those nodes. For every parameter combination $\mu^a$ and $\mu^r$ we performed 10 simulations in total. To quantify (i), we categorize opinion patterns into three categories: full consensus, bi-polarization, and non-trivial. For this purpose, opinion values at the end of the simulation period are apportioned by creating three equal-sized bins ([0, 1/3], [1/3, 2/3] and [2/3, 1]) covering the whole opinion interval [0, 1]. If all opinion values land in the center bin, we classify the opinion pattern as a case of "full-consensus". If all opinion values land in off-center bins, we classify the opinion pattern as "bi-polarization". Every other pattern is classified as non-trivial. To quantify (ii), we calculate the absolute error between the simulated opinion histogram of eleven equal-sized bins and the empirical opinion histogram of Europe (see Fig. 1).

This analysis is done for three different network structures in terms of average clustering coefficient $C$ and average shortest path length $L$: a Watts-Strogatz network with rewiring probability $p = 0.0001$ (high $C$ and high $L$), with rewiring probability $p = 0.01$ (high $C$ and low $L$), and with rewiring rewiring probability $p = 1$ (low $C$ and low $L$). Networks are generated using NetworkX [7].

## 3   Results

Figure 5 shows example opinion dynamics within the parameter space of non-trivial patterns, $\mu^a = 0.005$ and $\mu^r = 0.015$, produced on a Watts-Strogatz network with high average clustering and low average shortest path lengths. The balancing between the assimilative and repulsive mechanism in combination with a highly clustered network structure leads to opinion distributions in line with the empirical data shown in Fig. 1. The simulated patterns include a strong dominant opinion center, left and right extremal peaks, and off-center clusters. Figure 6 depicts the simulation space of non-trivial patterns, i.e. patterns other than full-consensus or bi-polarization, and the model fit with the empirical left-right political landscape of Europe, for three different network structures in terms of average clustering and average shortest path length. For a fixed confidence level $\epsilon = 1$ and critical level $\omega = 2$, non-trivial patterns are observed when the ratio between assimilative and repulsive rate of convergence is approximately 1/3. As network structure is changed from $p = 0$ (high $C$, high $L$) to $p = 1$ (high $C$, low $L$), to $p = 1$ (low $C$, low $L$), the space of non-trivial patterns becomes narrower and the model fit with empirical data becomes worse. We find the best model fits on network structures with high clustering, while average shortest paths seems to have less of an impact.

**Fig. 5** Multi-modal opinion dynamics. Shows example opinion dynamic produced for the combined social influence model on a Watts-Strogatz network with high clustering and low average path length ($p = 0.01$ and 200 nodes). The simulation setup is the same in all three examples: $\epsilon = 1$, $\omega = 2$, $\mu^a = 0.005$, $\mu^r = 0.015$

## 4 Discussion

In this study we extend the existing knowledge on social influence mechanisms [5] in networked systems by analyzing a simple model that can generate complex opinion distributions. For this purpose, we investigated a networked model of social influence that includes both assimilative and repulsive dynamics. We find that the model produces non-trivial patterns, i.e. opinion distributions other than full consensus or bi-polarization, when the forces of assimilative and repulsive influence are balanced and the underlying network structure is highly clustered. Simulations performed in this setting produce opinion distributions that share the multi-modal characteristics of empirical left-right political opinion landscapes observed in data collected by the European Social Survey program [4]. Generally, consensus and bi-polarization, as explained by traditional models of social influence, describe many social situations very well. For example, given a number of different communication channels all equally suitable, people choose a specific one and influence others around them to choose the same. Over time, consensus is reached and everyone in the group uses the same channel. The alternative to the above scenario would be the division into polarized groups. An extreme case of this is bi-polarization, where the groups find themselves on opposing ends of the opinion spectrum. Both outcomes, consensus and bi-polarization, can be considered stable outcomes: If social influence makes individuals more alike in their opinions, then consensus becomes stronger over time. Similarly, if social influence makes individuals more unlike in their opinions, then the opposing groups drift apart over time. These dynamics have been studied with great success and opened the door to studies of more complicated or as we call it, "non-trivial", opinion distributions. Given our results, we suspect such non-trivial opinion

(a) Watts-Strogatz graph with p = 0



(b) Watts-Strogatz graph with p = 0.01



(c) Watts-Strogatz graph with p = 1



**Fig. 6** Patterns observed on Watts-Strogatz graphs (200 nodes) with varying clustering coefficients $C$ and average shortest path lengths $L$. Left column shows patterns other than full consensus and bi-polarization produced. Right column shows model fit with the empirical left-right opinion landscape in Europe (see Fig. 1). Top row shows simulations on a Watts-Strogatz with $p = 0$, middle row with $p = 0.01$, and bottom row with $p = 1$ Opinion values are taken at $t = 200$. The confidence level and critical level are fixed at $\epsilon = 1$ and $\omega = 2$

distributions, as seen in empirical left-right political opinion landscapes (see Fig. 1), to be unstable outcomes by nature. However, given a certain balance of assimilative and repulsive forces and a clustered network structure, the pattern can be sustained for a prolonged period of time. In this respect, the multi-modal opinion distribution simulated with the combined model can be explained against the backdrop of the different opinion convergence speeds of the mechanisms. A right balance between the assimilative and the repulsive mechanism allows opinions to initially assimilate in almost all parts of the networks. If this assimilative dynamic is fully dominant, meaning that some parts of the network build polarized node pairs, the repulsive influence of these pairs can very slowly shift the network towards polarization. If, however, the assimilative force is too high, convergence dominates the occasional repulsive dynamics by the repulsive mechanisms and repulsive node pairs cannot develop. With parameter values as shown in Fig. 5 the majority of the network initially assimilates except in places where nodes start with highly opposing opinions and the repulsive influence is strong enough to locally exceed assimilation. Because the repulsive mechanism requires truncation, the repulsive dynamics of just a few polarized node pairs can slowly polarize the whole network. In this regard, network structure determines how quickly the repulsive dynamics can propagate through the network. Highly clustered networks may, thus, slow down this dynamic significantly. Finally, it has to be discussed whether the implemented mechanisms at the micro-level and the meso-structural conditions represent realistic assumptions. There is plenty of evidence from empirical studies to justify both assimilative and repulsive influence. The assumption of including both influences in a single model could be justified by a multiple arguments. For one, individuals are inconsistent, that is, individuals sometimes act in an assimilative way and sometimes in a repulsive way, depending on the topic or whom they interact with. Alternatively, individuals in a population could be heterogeneous, with some reacting more assimilative and others more repulsive. Lastly, it could be the case that the opinion distribution observed are the product of multiple sub-opinions, some of which are influenced more by assimilative force, and others more by repulsive force. Note that none of these possibilities exclude the others and all may apply. However, it is important to point out that the experimental evidence for repulsive influence suggests that it happens primarily in specific conditions. Our model assumptions take this into consideration, since assimilative influence is always present while repulsive influence is situational. Regarding the meso-structural conditions, we aimed at finding a compromise between complexity and realistic depictions of real networks. The Watts-Strogatz model allows us to create a realistic network structures in terms of average clustering and shortest average path lengths, with the general premise that real social networks also exhibit high average clustering and a low average shortest path length. The network structures generated by us have comparable properties in this regard, e.g. networks for the collaboration of scientists [3].

In conclusion, we suggest that multi-modal opinion distributions can be best understood as the result of opposing forces acting on agents in a network. A strong imbalance of these forces leads to either consensus or bi-polarization. However, if the opposing forces are in a certain balance and the underlying social network structure is highly clustered, then interesting and non-trivial opinion distributions can sustain themselves over some time.

# References

1. Akers, R., Krohn, M., Lanza-Kaduce, L., Radosevich, M.: Social learning and deviant behavior: a specific test of a general theory. In: Contemporary Masters in Criminology, pp. 187–214 (1995)
2. Asch, S.: Studies of independence and conformity: I. A minority of one against a unanimous majority. Psychol. Monogr. Gen. Appl. **70**, 1 (1956)
3. Barabási, A.: Network science. Philos. Trans. Roy. Soc. A: Math. Phys. Eng. Sci. **371**, 20120375 (2013)
4. European Social Survey ERIC European Social Survey (ESS), Round 9—2018. NSD—Norwegian Centre for Research Data (2019)
5. Flache, A., Mäs, M., Feliciani, T., Chattoe-Brown, E., Deffuant, G., Huet, S., Lorenz, J.: Models of social influence: towards the next frontiers. J. Artif. Soc. Soc. Simul. **20** (2017)
6. Gómez-Serrano, J., Graham, C., Le Boudec, J.: The bounded confidence model of opinion dynamics. Math. Models Methods Appl. Sci. **22**, 1150007 (2012)
7. Hagberg, A., Swart, P., Chult, D.S.: Exploring network structure, dynamics, and function using NetworkX. Los Alamos National Lab (LANL), Los Alamos, NM, United States (2008). https://networkx.org/
8. Heider, F.: Attitudes and cognitive organization. J. Psychol. **21**, 107–112 (1946)
9. Lorenz, J.: Modeling the evolution of ideological landscapes through opinion dynamics. Adv. Soc. Simul. **2015**, 255–266 (2017)
10. Macy, M., Kitts, J., Flache, A., Benard, S.: Polarization in dynamic networks: a Hopfield model of emergent structure (2003)
11. Mason, W., Conrey, F., Smith, E.: Situating social influence processes: dynamic, multidirectional flows of influence within social networks. Pers. Soc. Psychol. Rev. **11**, 279–300 (2007)
12. Myers, D.: Polarizing effects of social interaction. Group Decis. Making **125**, 137–138 (1982)
13. Rashotte, L.: Social influence. In: The Blackwell Encyclopedia of Sociology (2007)
14. Salzarulo, L.: A continuous opinion dynamics model based on the principle of meta-contrast. J. Artif. Soc. Soc. Simul. **9**, 13 (2006). https://www.jasss.org/9/1/13.html
15. Swire-Thompson, B., DeGutis, J., Lazer, D.: Searching for the backfire effect: measurement and design considerations. J. Appl. Res. Mem. Cogn. **9**, 286–299 (2020)
16. Sunstein, C.: Conformity. In: Conformity (2019)
17. Van Houtte, M.: School type and academic culture: evidence for the differentiation-polarization theory. J. Curriculum Stud. **38**, 273–292 (2006)
18. Watts, D., Strogatz, S.: Collective dynamics of 'small-world' networks. Nature **393**, 440–442 (1998)

# Observations on Modeling Social Identity: Suggestions to Address the Challenges of Social Identity

**Norman L. Johnson** (ORCID)

**Abstract** In the last two decades, social identity (SI) modeling and simulation have significantly advanced. They are building on and, in many cases, improving the over a half-century of validated SI experimental studies and theories. In this paper, observations on modeling and simulation of SI explore niches of additional opportunities based upon multiple perspectives: the evolution of social organisms, non-competitive theories of evolution, emergent properties of collective problem solving, advances in non-social computational modeling, epidemiological simulations, and complexity science. Based on these observations, specific recommendations are provided for expanding SI modeling and simulation. The main recommendation is to develop a general model of SI based on the observation that all social organisms share common traits, such as the innate drive to form SI or how individual states of uncertainty or stress trigger SI, but also recognize that complex species present more complex expressions of SI. Other recommendations are: SI models must accommodate that not all expressed SI traits have origins in or require higher fitness, all or many SI traits have triggers and maybe trigger thresholds that must be modeled, the inclusion of emergent group performance that may change SI behavior and strategies, and the development of a SI community model for research and realistic applications.

**Keywords** Social identity · Agent-based model · ABM · Diversity · Group performance · Emergent properties · Multilevel system · Evolution theories · Conformity · Complexity science

## 1 Introduction to Broadening the Approach to SI Modeling

All social organisms, almost by definition, can be said to express SI. Yet, there appear to be few researchers attempting to model what is common to all social organisms, particularly SI. The author's realization of the universality of social identity in social

N. L. Johnson (✉)
Referentia Systems, Honolulu, HI 96819, USA
e-mail: norman@santafe.edu

organisms became apparent while attending the 2002 *Self-Organisation and Evolution of Social Behaviour Workshop* [11]. By construction, the organizers included an equal representation of experimentalists and theoreticians/modelers. The publication of the workshop proceedings in 2005 captured why a unified approach to social organisms is beneficial: *"Self-organisation of social systems can be observed at all levels of biological complexity, from cells to organisms and communities. Although individuals are governed by simple rules, their interactions with each other and their environment leads to complex patterns. … The study of social systems from the perspective of complexity science leads to unusual results that show that, by self-organisation, complex patterns of behaviour may arise from very simple behaviour. By building these rules into certain computer models we develop a new type of understanding. This method may be applied to social systems of all kinds and of all organisms. Yet, so far, it has rarely been used among biologists. Moreover, biologists are little aware of the use of this method in the study of social systems in humans"* [11]. Much has changed since the writing of this introduction: agent-based models (ABM) in social sciences and SI modeling are common [18, 22]. Yet, while many citations of the workshop publication appear after 2005, none address a unified modeling approach to social organisms, even in biology. Notably, none seem to associate SI with social organisms, except humans.

While the text above argues for the use of complexity science, a multilevel and evolutionary analysis for modeling SI also has benefits and is captured in personal communication by J.J. van Bavel in 2018, *"I follow the logic of consilience laid out by E.O. Wilson, which is that a theory that operates successfully at multiple levels of analysis is more likely to be true and stand the test of time. On those grounds, I think there is a lot to be gained by not only looking at social psychological aspects of identity but seeing how these unfold at higher levels of analysis (social systems) and lower levels of analysis (the brain and cognition)…Moving up and down levels of analysis can generate new predictions and insights that might be hard to see if we always stick at the same level of analysis"* [4]. This quote adds an evolutionary and multilevel perspective to the discussion of SI modeling.

This paper examines SI modeling from various perspectives, including how the evolution of social organisms of different species represents different levels of adaptation of SI, matching the complexity of their environment. It proposes alternative approaches or missed opportunities for SI modeling and simulation. The goal is not to subsume the more than five decades of SI clinical experiments and theory but to explore the niches which may have been overlooked, mainly as they might be applied to mature ABM applications to advise and solve real-world challenges.

## 2   Observations on the Modeling Social Identity

The following section captures observations on the modeling of SI, followed by a section selection of publications that illustrate the observation and recommendations.

### 2.1   SI is Fundamental to all Social Organisms, Not Just Humans

A complete treatment of the multilevel evolution of SI is beyond the scope of this paper, but an observation supports the utility of such an exercise. An example from the *Social Behaviour Workshop* cited in the Introduction is the observation that all social organisms—from slime molds to social insects to social spiders to social mammals to lower and upper primates—exhibit a type of social copying when stressed or uncertain, capturing the transition from an individual activity to collective coordination. For example, when stressed from lack of water or nutrients, a slime mold (the social amoeba) shifts from independent behavior to coordinated action, including self-sacrifice—the extreme expression of SI, leading to propagation [7, 19]. On the other end of the evolutionary spectrum, humans are also observed to switch to social copying when uncertain or stressed [5, 26, 29]. This observation is revisited in the discussion of the CONSUMAT model in §2.3. Hence, the behavior of copying or imitating peers under uncertainty and stress is a candidate for a universal feature of SI in social organisms.

What if SI modeling started with the goal of capturing what is shared across all social organisms as a foundation on which to build more complex descriptions that are species-specific? This modeling approach is standard in the hard sciences, where dynamical theories (governing equations) are developed in the broadest descriptions, such as the equations of motion, followed by applying specialized constitutive models and simplifying assumptions to model specific problem areas. With the accessibility of extreme computing resources, the hard sciences have had even greater success in realistic modeling across many fields where simplified models combined with high-resolution simulations proved to be as good or better than complex models at a lower spatial resolution, e.g., ABM in epidemiology [8, 9], simplified constitutive models in continuum mechanics, and the direct numerical simulation of turbulence. A similar understanding is developing in modeling social behavior, particularly with ABM's advantages of self-organizing functionality [14, 18]. Is an opportunity being missed in SI modeling where more realistic SI behavior can be captured from simplified SI models combined with realistic, dynamic social networks generated or changed by the SI model? Recent publications and reports that address this question are provided in §3.3.

## *2.2  Behavior-Changing Social Identity Can Form from Trivial Differences*

Another aspect of simplifying SI models may involve considering that some aspects of an individual's SI are less complex and more flexible than is often argued for humans. How would this observation affect SI modeling? The unlikely answer may be found in evolutionary theory.

One misconception about the origins of social behaviors in primitive social organisms is that the details of their expressions are genetically pre-programmed. But a researcher of social wasps, Gadagkar, concluded from decades of research that ecological, physiological, and demographic factors dominate the influences of genetic relatedness in selecting for or against social traits [7]. This suggests that the expressions of SI may be more fluid than previously believed, even in the least complex social organisms. To generalize Gadagkar's conclusion: SI is an innate drive in all social organisms, but where the expressions depend on the species' complexity and local environment.

Many experiments show how humans can form strong and behavior-changing SI from minimal differences, such as experiments with children using random, trivial differences [24]. Akerlof and Kranton's 2000 paper summarizes: *"… competition is not necessary for group identification, and even the most minimal group assignment can affect behavior. 'Groups' form by nothing more than random assignment of subjects to labels, such as even or odd. Subjects are more likely to give rewards to those with the same label than to those with other labels, even when choices are anonymous and have no impact on [their] own payoffs. Subjects also have higher opinions of members of their own group"* [2]. Does the ease of formation of behavior-modifying SI from random attributes change SI models, and in what way?

This suggestion to modify SI models does not reduce the significance of over a half-century of experimental research on SI, particularly by the influential work of Tajfel [25] and the extensions of Tajfel's Social Identity Theory (SIT) after his early death, both in maturation and to the successful application of SIT in a variety of unrelated fields, as reviewed by Brown [4]. One resolution of the mature and validated SI theories with the above observations is that the behavior-changing formation of SI is an innate drive or need in all social organisms, but where the expression of the need is dependent on the social sophistication and environment of the species. One could also argue that Tajfel's SIT applies in social situations where a mature expression of SI is preexisting or the experimental design stimulates the strong formation of SI. But in experiments where random associations without payoff lead to SI formation, the innate need for the development of SI is triggered without recourse to competitive motivations.

Another argument by analogy on the possibility that SI is an innate need that finds a variety of expressions comes from the history of evolutionary theory. A common popular belief, even bias, is that all "evolutionary" features expressed in the animal species have an evolutionary significance of higher fitness during formation. Many academic papers have written justifications for an observed feature in a species simply

because of the assumption that if it occurred, there must be an increase in fitness due to the feature from selective evolution. A more mature evolutionary theory proposes that once the engine of diversity creation exists, the engine continues to create lasting diversity, even if the evolutionary selection pressure is lessened or absent [15]. Hence, the observed diversity in mature expressions of nature isn't exclusively because of evolutionary fitness but also because of the lack of evolutionary fitness and selection. For example, the extreme diversity of coloration in birds may not be associated with any increased fitness due to the coloration, but because the diversity creation of colors is *not selected* by an increase in fitness, and the diversity production engine freezes in evolutionary color changes.

When the above argument is applied to SI, possibly the innate need to form SI without payoffs or changes in self-esteem can induce behavior-modifying SI from trivial, non-competitive, random features. A possible characterization of this process is that the need for SI is an innate attractor in the individual, in complexity parlance, which requires a minimal stimulus to cause SI formation and where the expression of SI depends on the individual's internal state and external environment. There is nothing specific to human SI in this speculation. Hence the viewpoint provides a unified SI foundation for all social organisms. This innate SI attractor may have been overlooked as a universal, cross-species trait due to experimental designs that trigger mature expressions of SI. This observation leads to the next topic of triggers and thresholds in SI dynamics.

## 2.3 Triggers, Thresholds, and Habitual Behavior in SI Dynamics

There are unasked questions concerning experiments where SI occurs from minimal or random differences discussed in §2.2. What are circumstances in which a new SI is induced, or a pre-existing SI expression is triggered? Or when multiple SIs exist in an individual, what circumstances cause the expression of one SI over another? Or, more generally, what are the endogenous (individual) and exogenous (environmental) conditions that form or stimulate the expression of a SI or selection of one SI from multiple SIs? Is the formation of SI a gradual or abrupt process? Can the expression of a SI be a habitual state, not requiring rational choice? These questions become more relevant as the expression and management of multiple SIs within an individual are recognized and modeled [20]. While answering all of these questions is beyond this paper, this subsection examines the importance of modeling triggers and thresholds of SI behaviors and distinguishing between modeling conscious and habitual states.

An example of an ABM that best explores these questions was developed for consumer dynamics by Jager et al. in 2000 to implement a composite model from the many validated but niche behavioral theories [13]. The CONSUMAT model used three dominant behavioral models for individual choice: (1) bounded rationality, (2) social awareness and imitation of other consumers (peers), and (3) a rest state of

habitual behavior—the thoughtless repetition of prior choices. The CONSUMAT model was tested using an ABM on different social networks. As a weak form of validation but a significant achievement, the full spectrum of consumer buying dynamics is replicated with different parameter selections: high volatility in product choice, prolonged time volatility with instabilities, highly stable choices with a high diversity of product selection, and highly stable with low diversity of product selection.

A trigger of an individual behavioral state is implemented to initiate a specific decision process. The two triggers in CONSUMAT that initiate the individual's transition from habitual behavior to an activated decision state are (1) increased stress and uncertainty, leading to social imitation and copying (as discussed in §2.1), and (2) dissatisfaction from a historical comparison of needs fulfillment, leading to a rational choice of different options based on bounded awareness. These modeling choices capture the realistic behaviors: (1) an individual will sustain habitual behavior unless triggered to a heightened state of internal or external awareness, and (2) triggers of different internal states induce different types of behaviors.

Perhaps, one reason that habitual SI behavior appears not to be included in experimental studies is that the experimental design often induces an activated SI state, either consciously or unconsciously. The absence of SI habitual states in experiments appears to be carried over to the simulations of SI, as captured in the pre-review of the current state of SI models [22]. The above modeling observations can be applied to SI models: (1) some aspect of habitual behavior needs to be included, and (2) different triggers select between different types of behavior, including SI and non-SI behaviors.

In private communications with Jager, he shared that adding thresholds was necessary for the dynamic realism in CONSUMAT, where a threshold of a trigger captures a tipping point from habitual to behavioral change: a behavior does not gradually appear with a non-zero stimulus trigger, such as uncertainty, but first appears at a threshold level. Again, specific to SI models, what are the different SI behaviors and their triggers, and do they require a threshold before the behavior is expressed?

To provide a perspective on the above observations, a comprehensive framework for mapping and comparing behavioral theories in models of social-ecological systems was proposed in 2017 [21]. The framework is intended for applications in natural resource management, but the social-psychological framework proposed generally applies and shares goals and features of the CONSUMAT development from 17 years earlier. While the presentation does not include the concept of social identity—"identity" is only stated once in a long list of individual need states where *"Needs are motivational goals/factors for behaviour,"* social norms are cited as a crucial element of a person's behavior and central to social science disciplines. Overall, one main recommendation of the study is the necessity for a comprehensive model to switch appropriately between different behavioral modes, including habitual behavior.

While no mention of triggers appears in the framework paper, the one threshold reference is *"What defines a loss versus a gain is a threshold, or more precisely, a reference point that is a reflection of people's expectations or beliefs about past*

*outcomes."* An example of a habitual fisher agent provides an informative description, illustrating that threshold levels need not be fixed: *"Every time step that it brings back a catch and its needs are satisfied the behaviour becomes stronger and the threshold to switch to a different behaviour becomes higher. If the satisfaction drops below a threshold, the agent will start deliberating about alternative behaviour."*

In summary, a comprehensive SI model needs to have a rest state of habitual behavior as a foundation, with activated states of behavior with corresponding triggers and thresholds, based on internal states and external influences.

## 2.4 Emergence and Emergent Properties in SI Group Utilities

The word "emergence" has become a common descriptor in many social science publications; for example, in 2008, *"90% of papers on complexity and social simulation explicitly refer to emergence"* [23]. Emergence is now commonly used to mean appearance, expression, coordination, and, possibly the least useful, surprise, and consequently has lost its technical meaning [3]. This widespread usage of emergence does not capture the definition for an emergent multilevel property: a feature observed in the group (global) but not observed or expressed in the individual (local).

For most modeling studies of SI, the goal is to provide a descriptive model of known or proposed SI features for evaluation, where the expression of SI or its utility is not treated as an emergent property. One example of a limitation of not including emergent dynamics and features in the modeling is when the utility of the group has an emergent component but is not captured, which can, in turn, cause the lack of the individual utility to reflect the full expression of the group utility and, therefore, might change the conclusions of the study. This limitation is in addition to the additional difficulty that if the group utility is explicitly modeled within the individual, the question arises as to the realism of the modeling: the group utility cannot typically be objectively known because individuals have only perceptions of group utility but no mechanism to evaluate the group state objectively. The exception to this statement is when group payoffs are explicitly made to individuals by an intentional group structure.

An excellent example of the hazards of omitting emergent properties is the decades of studies on the evolutionary origin of cooperation in publications. Many of these studies largely fail in their goal because the models explicitly include cooperative behavior as an option within the individual behavior. In this explicit modeling approach, the simulations cannot demonstrate the emergent origin of cooperation but only the desirability or selection of cooperation. By contrast, if an agent behavioral model doesn't include cooperation, but the global dynamics of the simulation exhibit emergent cooperation, then the model and simulation can be strongly stated to capture the origin of emergent cooperation. Then, by using evolutionary processes, once the emergent property increases individual fitness, the emergent cooperative expression can be internalized within the population of individuals through selective genetics [15].

The ABM simulations of Hemelrijk in 1997 of the dynamics of herd structure [10] illustrate the above argument. In simulations with only aggressive individual behavior, Hemelrijk observed that a stable interaction could occur between a strong individual and multiple weaker individuals in the formation of the dominance structure of the herd. The multiple weaker individuals exhibited emergent cooperation, even though the behavioral model did not include individual cooperative behavior. Many models at the time claimed to demonstrate that cooperation was an evolutionary adaptation to higher fitness. Yet, the individual models typically included cooperation as an individual option and arguably failed in their demonstration [10].

While the evolutionary origins of SI may be less attractive to many researchers, the above discussion has relevance to SI modeling choices and possibly SI theories. For perspective, one of the significant advancements in evolutionary theory in the last two decades is the resolution of the controversy concerning group utility in evolution, as captured in a monograph by two of the most influential evolutionary theorists, Wilson and Wilson, in 2007: *"Current sociobiology is in theoretical disarray, with a diversity of frameworks that are poorly related to each other. Part of the problem is a reluctance to revisit the pivotal events that took place during the 1960s, including the rejection of group selection and the development of alternative theoretical frameworks to explain the evolution of cooperative and altruistic behaviors… Multilevel selection theory (including group selection) provides an elegant theoretical foundation for sociobiology in the future, once its turbulent past is appropriately understood"* [28]. Although SI should be a key component of sociobiology theories, it is not mentioned in the monograph. While this omission is significant to the history of SI theories, the discussion of its implications is beyond the scope of this paper. Still, specific aspects of the multilevel evaluation of utilities are relevant to ABM SI modeling and can be discussed.

The key to determining utilities in the context of SI is capturing the benefits and costs expressed at multiple levels: for agents, an SI group of agents, and communities of SI groups. A feature of all ABM treatments of SI is the use of agent and group utilities, either as payoffs or for strategy evaluations. For example, if a rational choice model is used, then the utility of an agent determines the agent's behavior. If different individual or group management strategies are examined, group utilities are used to evaluate them. While it is beyond the goal of this paper to review the models of SI utilities, such as the commonly used self-esteem [25], the emergent sources of utility appear to be overlooked in multilevel SI models of individual and group(s).

In the late 1990s, two groups of researchers independently discovered how diverse groups could outperform the average individual and how even a group of high-performing experts can be outperformed by a group of individuals with a diversity of individual performance or skills [12, 14]. Identical to the challenges faced in gaining acceptance of group selection described by Wilson and Wilson [28], both first attempts to publish these results were rejected, with a reviewer of my 1998 submission stating, *"I don't see what is wrong, but it can't be right."* Two decades later, these concepts are popularly accepted and published as "collective intelligence"

and are key to understanding the invisible hand in optimizing stock markets and managing large research programs [16]. The following asks if a similar bias has occurred in the history of SI modeling.

One example of emergent utility is when an optimal but emergent group solution to a problem may not be comprehensible to the individual. In a 1998 report, I analyzed how information derived from a collection of independent agents solving a maze can be aggregated to obtain the shortest path [14]. Because a myotic agent has no global perception of the maze, the agent has no mechanism to judge the quality of its chosen path. A significant discovery was that any reduction in the contribution of experiences by the agents in the aggregation for the group solution led to reduced group performance. This discovery led to an analysis that found that group performance correlated with the diversity of individual contributions to the group solution. This diversity correlation occurs only for a range of problem complexity that confounds an expert solution but is not so great as to cause the individual's contribution to be noise [16].

In the 1998 study, it was assumed that all agents had a common worldview (they agreed on options in the maze), reflecting a common SI. In a later study of the same maze problem but using agents with different worldviews or SI (they disagree on options), the resulting biases lowered the group performance unless the biases themselves were diverse or, more accurately, uncorrelated [16]. An additional discovery of the 1998 study was that the optimal emergent group performance was when each individual could communicate their full experience to the group solution, not their best option, nor a uniform weighting of all options. One way to understand these results is that in complex problem domains, individuals have diverse and non-overlapping areas of experience. One individual, including an expert, cannot perceive the global problem in complex problem domains. The collective aggregate of experience or skills always yields a better solution than an average performer and often the expert.

These results have direct application to the SI modeling: (1) emergent group utility can be uncorrelated with aggregate individual utility, which in turn, may alter conclusions about the efficacy of SI, (2) a higher emergent utility of a SI group requires compatibility of individual contributions—a shared worldview or SI, (3) because the emergent solution is robust to uncorrelated bias and even extreme noise in the individual contributions [14], SI groups may show higher emergent performance in experiments in the presence of miscommunication, misinformation, or low SI coherence, and (4) optimal SI group performance occurs when individuals of a SI group can communicate their complete experience, which could be restricted by repressive SI conformity. In summary, including emergent properties in multilevel SI simulations can result in more robust and realistic models, change the conclusions of studies, and contribute a new understanding of SI in group performance.

## 3   Illustration of the Above Observations to ABM SI Studies

This section examines three recent papers describing ABM implementations of SI theories to illustrate the observations of the prior sections. These studies were selected based on the quality of the behavioral models and implementation choices, representing this author's view of the sophisticated state of SI modeling. While few papers were selected to illustrate the observations presented herein, the advantages of the observations are hopefully helpful to other publications and identify SI modeling additions for more realistic applications.

### 3.1   ABM of a Comprehensive Social Identity Theory (SIT)

Upal and Gibbon, in 2015 [27], presented a socio-cognitive model of SI dynamics and illustrated how agent-based social simulation could be a valuable tool for theory refinement. The simulations use a rational choice theory that maximizes individual utility. Intergroup behavior is driven by the need to maintain positive self-esteem, derived partially from affiliation with SI groups. Comments on the implemented SIT's accuracy are beyond this paper's scope, but the study is an example of the advanced implementation of a mature behavioral model. The SIT model captures a comprehensive spectrum of socio-structural beliefs, individual and collective strategies, intergroup permeability, and personal and group costs… to name some of the features. The simulations of 100 agents examined 12,000 simulation groups with 500 rounds per group, initializing each run with random distributions of individual resources, agent perceptions of permeability, legitimacy, stability, and individual esteem. The analysis of the simulations examined correlations between the input variables and outcomes of multiple SI management strategies. Given the maturity of the SIT model, the analysis provided extensive results on the sensitivity of different strategies to the model parameters. The strongest correlations observed were that out-group resources were negatively correlated with all SI management strategies. *"This means that agents are more likely to denigrate, glorify, attack and change entry conditions targeting groups that are believed to have few resources"* and *"As in-group resources increase, agents become more likely to engage in collective strategies against the out-group members."* The two unexpected results, labeled "emergent," were (1) the positive correlation between average group resources and all SI actions and (2) the negative correlation between outgroup resources and SI actions.

The reason for citing this study is to note that the implementation of the SIT is linear in all relationships (an explicit assumption) and deterministic (the same initial conditions produce the same outcome). The model excludes triggers and thresholds in behavior, which would introduce nonlinear dynamics. Similarly, there is no modeling of habitual behavior, which adds a strong path dependency in the solutions, another nonlinear behavior. The deterministic nature of the model excludes the possibility of SI forming from random events. The addition of modeling any of these behavioral

effects while increasing the complexity of the analysis would result in a more realistic model and results. A final comment is that the unexpected results are labeled emergent patterns, using the more popular definition of emergence. There is no indication in the results that the simulations show emergent behavior as defined in §2.4.

## 3.2 ABM Study of Trust and Conformity, Using Fitness of Group Diversity

A 2022 paper by Fazelpour and Steel studies the positive and negative effects of different types of diversity on SI performance using an ABM [6]. The problem challenging each agent is selecting two options with unknown payoffs that are sequentially observed to optimize their preference. Their resulting payoff preferences can be shared based on a predetermined and fixed social network. The study's main conclusions are that different types of diversity *"can, in certain circumstances, benefit collective performance by counteracting two types of conformity that can arise in homogeneous groups: those relating to group-based trust and those connected to normative expectations toward in-groups."* The main conclusions duplicate the earlier diversity studies described in §2.4 and [14]. Still, because the simulations include multilevel SI dynamics of information sharing and blocking, the nuances of the effect of diversity on collective SI performance are also revealed. While the use of a fixed social network does not realistically represent SI group formation and change, as discussed in §2.2, the authors' variable weights of social network connections are stated to capture intergroup dynamics, but no details are provided. No modeling information is provided if triggers and thresholds were included in implementing behavior models, communication, or strategies. Habitual behavior is not mentioned.

## 3.3 Multipurpose SI Community Model for Large-Scale ABM Simulation

A significant advancement of epidemiology and its usefulness in pandemic strategies transpired in the 2000s when ABM simulations with billions of agents were demonstrated at Los Alamos National Laboratory by modifying a molecular dynamics simulation resource. The resulting ABM epidemic modeling resource, EpiCast, simulated pandemics at a national level, capturing the movement and infection state of every individual in the U.S. (300 million at the time) using census and mobility data [8]. The EpiCast results were so influential that pandemic policy decisions of the last century were changed in the U.S. and internationally and have continued today with the COVID pandemic, utilizing the rapid development of vaccines instead of a national quarantine. A critical precursor that made EpiCast possible was developing

a 2000-person ABM community model of the infectious spread of smallpox [9]. The advantage of the community model is that it captures the realistic spread of infection through a contact network with movement between homes, workplaces, and public locations. The model was validated with other infectious diseases and became a standard test platform for developing new infectious models. EpiCast replicated this model to duplicate the populations of each county in the U.S., thereby capturing the entire U.S. population.

Based on the success of EpiCast as a team member and PI, I developed a research proposal in 2009 [17] after concluding a Phase 1 exploratory study for an ABM resource for managing message campaigns in actual geographical regions with polarized SI populations, using a simplified SI model, a replicated community model based on the smallpox community model [9], and data-driven social networks. The combined ABM resource with data assimilation was argued to assist decision-makers in conflict management and policy deployment. Another trial SI community model was proposed in 2022 to study the "emergence of social norms" [1]. This study also adds genetic algorithms to enable the evolution of rules to optimize individual fitness in the presence of information exchange, enabling the discovery, rather than a specification, of collective norms.

The dynamical similarity between a community experiencing an infectious disease with adaptive behavioral changes and a community experiencing SI formation and adaptive behavioral changes suggests that the development and use of a SI community model might be transformational to the testing of new SI models and the development of large-scale policy management resources, similar to the experience of EpiCast.

## 4   Conclusions and Future Studies

These are highlights of the suggestions that might be included in future SI resources. §2.1: Start with a universal SI model common to all social organisms, and then specialize the model for specific social organisms—the more complex the organism, the more complex the SI. §2.2: Consider that expressions of SI may not require modeling of fitness but can occur by chance, reflecting the attractor nature of SI. §2.3: Consider inclusion of habitual behavior and what triggers and thresholds activate each SI feature. §2.4: Allow for emergent properties in multilevel SI models, particularly in how group performance benefits the individual. And, §3.3: Consider the development of a validated, multi-purpose SI community model with realistic, highly-resolved, SI-driven social networks.

A theme throughout this paper is that the challenge of the high complexity of evolved human SI may hamper the advancement of SI modeling. And how an evolutionary perspective might guide the development of SI models. For this author, the most exciting discovery in examining the evolutionary development of SI is the perspective that human SI might be viewed as an emergent collective consciousness of the group. This observation aligns with an unpublished theory of the author that the evolutionary origin of consciousness or sentience in an organism is the ideation

equivalence of the biological sense-of-self of advanced immune systems to address the high internal complexity of a multicellular organism. From this viewpoint, SI evolved as an expression of emergent immunity of the SI group to outside ideas while managing the SI group's high internal complexity or diversity. This leads to the observation that in lower forms of social organisms, SI is not self-aware or emergent but purely responsive at an individual level. And, in higher social organisms, emergent SI provides forms of group awareness and immunity to outside ideas, which the individual cannot understand.

# References

1. Agrawal, R., Ajmeri, N., Singh, M.P.: Socially Intelligent Genetic Agents for the Emergence of Explicit Norms. ArXiv abs/2208.03789 (2022)
2. Akerlof, G., Kranton, R.: Economics and Identity. Q. J. Econ. **115**(3), 715–753 (2000)
3. Bedau, M.A., Humphreys, P.: Emergence: Contemporary Readings in Philosophy and Science. MIT Press, Boston (2008)
4. Brown, R.: The social identity approach: appraising the Tajfellian legacy. Br. J. Soc. Psychol. **59**(1), 1–21 (2019). https://doi.org/10.1111/bjso.12349
5. Cialdini, R.: Influence: Science and Practice. Allyn & Bacon, Needham Heights (2001)
6. Fazelpour, S., Steel, D.: Diversity, trust, and conformity: a simulation study. Philos. Sci. **89**(2), 209–231 (2022). https://doi.org/10.1017/psa.2021.25
7. Gadagkar, R.: The Social Biology of Ropalidia Marginata: Toward Understanding the Evolution of Eusociality. Harvard University Press, Cambridge (2001)
8. Germann, T.C., Kadau, K., Longini, I.M., Macken, C.M.: Mitigation strategies for pandemic influenza in the United States. Proc. Natl. Acad. Sci. U.S.A. **103**, 5935–5940 (2006)
9. Halloran, M.E., Longini, I.M., Jr., Nizam, A., Yang, Y.: Containing bioterrorist smallpox. Science **298**(5597), 1428–1432 (2002)
10. Hemelrijk, C.K.: Cooperation without genes, games or cognition. In: Harvey, P.H.A.I. (ed.) Fourth European Conference on Artificial Life. MIT Press, Cambridge (1997)
11. Hemelrijk, C.K. (ed.): Self-Organisation and Evolution of Biological and Social Systems. Cambridge University Press, Cambridge (2005)
12. Hong, L., Page, S.E.: Groups of diverse problem solvers can outperform groups of high-ability problem solvers. PNAS **101**(46), 16385–16389 (2004)
13. Jager, W., Janssen, M., De Vries, H., De Greef, J., Vlek, C.: Behaviour in commons dilemmas: implementing the CONSUMAT approach in an ecological economic model. Ecol. Econ. **35**, 357–379 (2000)
14. Johnson, N.L.: Collective Problem Solving: Functionality beyond the individual. Los Alamos National Laboratory Report, LA-UR-98-2227, Los Alamos, NM (Nov 1998) www.academia.edu/download/84598680/NLJsims_AB_v11.pdf
15. Johnson, N.L., Watkins, J.H.: Interplay of adaptive selection and synergistic performance: as an example of natural selection and self-organization. In: Selection and Self-Organization Workshop, CSIRO-Sponsored Complex-System Science Workshop. Katoomba, Australia (2007). https://doi.org/10.2139/ssrn.2232193
16. Johnson, N.L.: Applied science of collective intelligence: solving the grand challenges facing humanity. Spanda J. **5**(2), 97–108 (2014)
17. Johnson, N.L.: Final report of SAGE: situational awareness for the GTWO environment project, 98 pages (September 23, 2009) CDRL A001AC, N00014-09-M-0190, Office of Naval Research, 8725 N. Randolph, Arlington, VA 22203
18. Miller, J.H., Page, S.E.: Complex adaptive systems: an introduction to computational models of social life. In: Princeton Studies in Complexity. Princeton University Press, Princeton (2007)

19. Ostrowski, E.A.: Enforcing cooperation in the social amoebae. Curr. Biol. **29**(11), R474–R484 (2019)
20. Reed, A., II.: Special session summary: when what I think feel and do depends on who I am: identity effects on judgment, choice and self-reinforcement. Adv. Consum. Res. **31**, 335–338 (2022)
21. Schluter, M., Baeza, A., Dressler, G., Frank, K., Groeneveld, J., Jager, W., Janssen, M.A., McAllister, R.R.J., Mueller, B., Orach, K., Schwarz, N., Wijermans, N.: A framework for mapping and comparing behavioural theories in models of social-ecological systems. Ecol. Econ. **131**, 21–35 (2017) doi.org/https://doi.org/10.1016/j.ecolecon.2016.08.008
22. Scholz, G., Eberhard, T., Ostrowski, R., Wijermans, N.: Social identity in agent-based models: exploring the state of the art. In: Ahrweiler, P., Neumann, M. (eds.) Advances in Social Simulation ESSA 2019. Springer Proceedings in Complexity. Springer, Cham (2021)
23. Squazzoni, F.: Book review of emergence: contemporary readings in philosophy and science. JASSS **11**, 4 (2008). www.jasss.org/11/4/contents.html
24. Stangor, C., Jhangiani, R., Tarry, H.: Ingroup favoritism and prejudice. In: Principles of Social Psychology, 1st International H5P Edition (2022). opentextbc.ca/socialpsychology/
25. Tajfel, H., Turner, J.C.: An integrative theory of intergroup conflict. In: Austin, W.G., Worchel, S. (eds.) The Social Psychology of Intergroup Relations, pp. 33–47. Brooks/Cole, Monterey, CA (1979)
26. Tesser, A., Campbell, J., Mickler, S.: The role of social pressure, attention to the stimulus, and self-doubt in conformity. Euro. J. Soc. Psychol. **13**, 217–233 (1983)
27. Upal, M.A., Gibbon, S.: Agent-based system for simulating the dynamics of social identity beliefs. In: Proceedings of the 48th Annual Simulation Symposium, pp. 94–101 (2015)
28. Wilson, D.S., Wilson, E.O.: Rethinking the theoretical foundation of sociobiology. Q. Rev. Biol. **82**(4) (2007). https://doi.org/10.1086/522809
29. Wooten, D.B., Reed, A., II.: Informational influence and the ambiguity of product experience: order effects on the weighting of evidence. J. Consum. Psychol. **7**(1), 79–99 (1998)

# The Friendship Field - an Agent-Based Model on Dyadic Friendship Formation Driven by Social Battery

**Chrisja Naomi van de Kieft** and **Eva Margretha Timmer**

**Abstract** Humans have an intrinsic need for friendship, especially in adolescence when entering a new social environment where they do not know anybody. The question as to how friendships form is frequently asked. In research, three important factors have been identified in the formation of friendship: extraversion, resemblance and social status. To our best knowledge, a missing aspect in current research on friendship formation is the concept of "social battery". The social battery comprehends an individuals' energy level to engage in social contact. When the social battery is exhausted, it can prevent an individual from social contact, and consequently from making new friends. The recharging and exhaustion of the social battery heavily depends on the person's extraversion level. In this paper, we develop an agent-based model "the Friendship Field" that simulates real-life dyadic friendship formation where the individuals' interactions are motivated by their social battery. With this model, we investigate emergent patterns regarding extraversion, resemblance and status. The model reproduces a pattern of the mere-exposure-effect, an existing theory on friendship formation. Moreover, it proposes a new factor for friendship formation in social sciences: the social battery.

## 1  Introduction

Humans feel the need to belong, the need for closeness and thus, the need for friendships [1]. Especially in adolescence, the importance of friendships increases [3]. They provide a sense of acceptance and belonging for young people. Friendships

C. N. van de Kieft · E. M. Timmer (✉)
Department of Social Sciences, Information Technology Group, Wageningen University and Research, Hollandseweg 1, 6706 KN Wageningen, The Netherlands
e-mail: eva.m.timmer@hotmail.com

C. N. van de Kieft
e-mail: chrisjavdkieft@gmail.com

299

help to develop compassion, care, and empathy, as well as a helping to define the self, outside of the family. A very important year for adolescents to establish new friendships is in the first year of university. This is an example of entering a new social environment with many unknown people. The question of 'How do friend-ships form?' has been the topic of interest in a broad range of research fields for over 100 years [2, 5, 10]. In many of these studies, several social interaction aspects were nominated as having a significant impact on friendship formation. However, to our best knowledge, the individuals' energy level to engage in social contact is not considered in previous studies. We propose to implement this energy level as a "social battery". In this paper, we develop an agent-based model "the Friendship Field" that simulates real-life dyadic friendship formation where the individuals' interactions are motivated by their social battery. People can only become friends when they have enough energy for social contact.

With our model, we investigate emergent patterns regarding a set of well-studied factors in the field of social sciences: extraversion, similarity (or resemblance) and social status [11, 12, 14, 15]. Furthermore, the model proposes a new factor in friendship formation: the social battery. The Friendship Field simulates a generic version of friendship formation, since it can be applied to all situations with no initial friendships and no hierarchy. However, the model was built according to the idea of first year bachelor students in the Netherlands, therefore in this cultural setting status is of less importance than in other situations.

The remainder of this paper is organised as follows: Sect. 2 (Background) elaborates on the factors influencing friendship formation and introduces the concept of social battery. Section 3 (Model Design) provides an overview of the model design including the timestep description and parameters. Section 4 (Model Validation and Analysis) presents the results of the sensitivity analysis. Section 5 (Discussion and Conclusion) summarises the findings and proposes suggestions for further research.

## 2   Background

Extraversion is a personality trait in the "Big Five" OCEAN model of personality by Costa and McCrae [4]. The model describes a range from a passive and reserved personality to more talkative and sociable characteristics. People with low extraversion (introversion) prefer smaller groups or solitude and avoid large social situations, while people with high extraversion tend to seek the company of others. Extravert people experience social interaction more often as positive and are therefore more likely to select friends [15]. Another important difference between extraverts and introverts is that extraverts gain energy from social interaction, while for introverts it costs energy [8].

The tendency to become friends with people that are close in space to you is called the proximity principle [6]. Proximity (literal closeness) increases the interaction between people, which is key for friendship formation [10]. People with similar

characteristics (figurative closeness) have a higher chance of reciprocating a friendship nomination [12]. A friendship nomination is considering someone a "friend", so when the nomination is reciprocated a bidirectional friendship arises [11].

This phenomenon of bidirectionality is inherently intertwined with status. Status encompasses many aspects of social life, including social importance, appeal, and kindness [7]. In Kemper's status-power theory, status is the voluntary compliance with what the other human wants [9]. Kemper's status-power theory of social relationships describes (amongst others) the balance between gaining status, conferring status, and claiming status. Roughly said, when a person receives a status conferral, they compare the conferral to their expected and wanted status. If these match, the person accepts the status conferral. Two people who accept each other's status conferral consider the interaction as pleasant and therefore, their relationship is improved. When the status conferral is lower than expected, the recipient considers the interaction unpleasant. The power part of the status-power theory is based on involuntary compliance, when the expected status is not given so the human forces the interaction to gain status. However, we feel that this does not happen in our cultural environment with no hierarchy and therefore we will not consider this in our model.

People do not always have the energy to engage in social contact. The level of this energy can be regarded as the "social battery", since it can deplete and recharge. In grey literature, the term social battery is referred to as one's capacity to engage in social contact [13, 16]. For people with various levels of extraversion, the social battery depletes or recharges at a different pace during social interaction. The state of the social battery might influence whether one interacts with people and thus, whether friendships can arise. Since social contact is a key factor in friendship formation, we believe that social battery is a crucial component of the friendship dynamics. Moreover, interactions with known people cost less or bring more social energy than interactions with new people. Therefore, a lower social battery level can prevent an introvert person from interacting with yet unknown people but might allow for interactions with already established friends. Thus, the ways in which the social battery affects different people and their interactions, influences the friendship dynamics.

Furthermore, Kemper theorizes about reference groups. Reference groups are the groups of people who motivate us to perform actions [9]. Reference groups can be anything or anyone, alive, dead, real or imagined. Besides other people as a reference group, one of the reference groups mentioned by Kemper, is the "organism" or the body. The body has certain interests (food, energy, sex) and these interests can be regulated or repressed by social or moral codes. When the organism is pushed to its limits, it can refuse. We linked the concept of social battery to the organism as a reference group, where the limiting factor for engaging in social interaction is the organism's social energy.

Based on these concepts, friendship networks arise in our model. In these networks, people have a different number of friends and stronger or weaker relationships. As in real life, relationships can decay over time when a long time of no interaction occurs. With no repeated interaction, the relationships fade over time, and this gives the network a dynamic structure.

# 3 Model Design

The full model can be retrieved from https://www.comses.net/codebases/29cb34a1-9128-494d-937d-02d1e34b5fc4/releases/2.0.0/.

## 3.1 Timestep Description

**Setup**. At the setup of the model, the humans are placed on the field. Randomly, all humans receive an extraversion level ranging between −1 and 1, and a status, kindness and appeal ranging between 0 and 1. The humans receive a list of random characteristics that can either be a 0 or a 1 (characteristics-list). This number represents the absence or presence of a characteristic (e.g. are they vegetarian?) that the humans possess. Furthermore, their initial social batteries (social-battery) are set to the social battery threshold (sb-threshold), a slider in the interface. An extraversion below 0 means the human is an introvert and a positive extraversion means the human is an extravert. Extraversion levels in the range −0.1 to 0.1 are excluded to create more separation between the humans.

    **Chill on the field**. On the start of a timestep, the humans "decide" whether they want to meet someone (Fig. 1). Their want-to-meet is a probability ranging from 0 to 1, calculated with their social-battery, their extraversion level and the sb-threshold. For extraverts, a social-battery level higher than the sb-threshold means that they have a sufficient level of social energy and do not necessarily need to meet another human



**Fig. 1** Flowchart of one timestep in the friendship field

to involve in social contact. For introverts, a social-battery level higher than the sb-threshold means that they have enough social-battery to engage in social interaction. A social-battery that is below the sb-threshold means for extraverts that they need to have social contact, while in this situation the introverts would rather avoid social contact. If the humans do not engage in social interaction, they chill on the field (COF). Half of the human's extraversion level is added to the relative difference within the social-battery and the sb-threshold, as extraversion influences the desire to meet people in real life.

$$RelativeDifference = (1 - abs(SocialBatterySBThreshold) / (MaxSB - MinSB))/2$$

$$WantToMeet = RelativeDifference + Extraversion/2$$

(N.B. In Netlogo, different words within a variable name are separated with a dash '-' by convention. For readability, the variables in the formula boxes are in CamelCase style). The want-to-meet is bound to fall between 0 and 1. A human's not-want-to-meet equals 1—want-to-meet. The human boolean variable is-available? is set to True or False, weighted with the want-to-meet and not-want-to-meet probabilities.

**Finding a mate**. Next, the humans try to find a mate: another human with whom they can interact. The probability of meeting someone is related to the relationship level (relation). This means that the better the relationship between two humans is, the higher the chance they meet again. Each human semi-randomly picks a candidate from all available humans and the humans become each other's mates. The extravert humans immediately interact, while the introvert humans first check whether the mate is already a friend. If the mate is already a friend, the interaction starts. If the mate is not yet a friend, the introvert humans first check if they have enough social battery to engage in a social conversation with someone they do not know. This is dependent on the extraversion level of the human and the introversion-weight, another slider in the interface.

$$SocialBattery + Extraversion * IntroversionWeight >= SBThreshold$$

This social battery check is because meeting new people generally costs more social battery then already known people, especially when the human is introverted.

**Status conferral**. The first step of the interaction is the status conferral. The mate of the human determines how much status they are willing to confer to the human based on status, appeal of the human, kindness of the mate and existing relationship status between the human and the mate [7].

$$StatusConferral = ((StatusWeight * MatesStatus) \\ + (PersonalityWeight * (OwnKindness + MatesAppeal)/2) \\ + (CurrentRelationWithMate * RelationWeight))$$

The weight of the status (status-weight), personality (personality-weight) and relationship (relation-weight) can be changed with the sliders in the interface. These weights sum up to 1, but their relative weights can be adjusted to investigate the importance of the three determinants in the status conferral. If the status-conferral of the mate is higher than the status of the human, this will have a positive effect on the interaction. When the received status-conferral is lower than the human's status, the human is offended, and this has a negative effect on the interaction.

**Interaction**. After the status conferral, the interaction itself takes place. The human and its mate both determine the interaction value based on the percentage of corresponding characteristics, extraversion level combined with the relationship (openness), and the status-conferral. As you get to know people (relation increases), your extraversion contributes less to your openness. Extraverts have an advantage in the relationship increase when the relationship is yet beginning, while introverts strengthen their relationships with people they already know more easily. The similarity in characteristics and openness both have a higher weight in the equation than status conferral, according to the cultural setting of this model. However, this can be adjusted in the code.

$$Interaction = (Percentage * CharsWeight + Openness * OpennessWeight$$
$$+ StatusConfWeight * StatusConferral)/10$$

**Update**. After interaction, three variables of the humans get updated. The status is updated based on the status conferral they received from the mate. The volatility is hard-coded as 0.10, according to the start value in the playground model [7].

$$Status = (Status + (StatusConferral - Status) * Volatility)$$

The second variable that gets updated is the relation between the human and the mate. The new relationship is increased with the interaction value. When the relation exceeds the friendship-threshold, the humans are considered friends.

$$Relation = Relation + Interaction$$

Thirdly, the humans update their social battery, see section 'Social battery and Extraversion'.

**Un-mate**. In un-mate, the humans set their mate to nobody and the relationship-links that are above the friendship-threshold become visible.

**Relation-decay**. At the end of the timestep, the relationships of all humans decay. This formula takes into account the yet established relation value between humans, because relationships with people you are very good friends with have a lower decay rate.

$$Relation = (Relation - (1.1 - Relation) * 0.05)$$

The 1.1 prevents the decay from being 0. Whilst calibrating, the 0.05 is a decay-constant that provided the most plausible outcomes in our opinion.

**Social battery and extraversion**. For introverts, it costs energy to have social contact and the social battery recharges on the field. On the contrary, extraverts gain energy by engaging in social contact and exhaust their social battery whilst on the field. We implemented this in our model by creating two different ways for the humans to lose and gain social battery. The social battery changes after the interaction and when the humans chill on the field due to not having found a mate.

$$Field \, Social \, Battery:$$
$$Social \, Battery = Social \, Battery - Field \, Battery \, Constant * Extraversion$$
$$Interact \, Social \, Battery:$$
$$Social \, Battery = Social \, Battery + Interact \, Battery \, Constant * Extraversion$$

The interact-social-battery ensures that the extravert humans gain energy from interacting, and the introvert humans lose energy from interacting. This same mechanism is used the other way around in the field-social-battery.

**Adaptive characteristics**. When interacting with friends with different characteristics, people tend to adopt each other's characteristics. To implement this, a switch adaptive-characteristics? was built which turns on the adaptation of characteristics. At the end of the interaction, the human checks whether the mate is a friend. If this is true, the human with the lowest status changes one characteristic to its mate's.

**Run time**. One timestep resembles one chance to interact with someone. We assume three interaction chances per day, therefore one day consists of three timesteps. As the study investigated how friendships formed during one year, the model runs for 1095 timesteps.

## 3.2   Parameters, Variables and Design

**Agents**. The agents in this model are humans. These humans have a fixed extraversion level, kindness, appeal and a list of abstract binary characteristics. They also have an adaptive status, social battery and current mate with whom they interact. Besides this, the humans also have a variable candidate. This variable is used to determine the other humans that are available to become a mate with. Based on these interactions, the humans develop a certain relation and hence, the humans can form friendship networks. All human variables are described in Table 1.

**Environment and Networks**. The grid is a theoretical space in which humans can meet but it has no spatial meaning. The humans are able to confer status to each other and depending on their interaction, their mutual relationship can change. All the humans have a link to all the other humans and the humans function as nodes, their reciprocal relationship represented by the weight of that edge. Furthermore, the links all have a resemblance: the percentage of agreeing characteristics they possess. The

**Table 1**  Overview of the variables attributed to the humans

| Variable | Use | Value |
|---|---|---|
| Extraversion | Level of how extravert a human is | −1 to 1 |
| Status | Amount of status a human has | 0–1 |
| Kindness | How easily they confer status to others | 0–1 |
| Appeal | How easily they attract status from others | 0–1 |
| Social-battery | How much social energy they have to form new friends/maintain friendships | 0–100 (START = 100) |
| Characteristics-list | List of binary characteristics in which they can resemble others | List of random 0 s or 1 s |
| Candidate | Possible human to become mates with | Another human (START = nobody) |
| Mate | The human with whom they are interacting | Another human (START = nobody) |
| Is-available? | Whether the human is available to meet | TRUE or FALSE (START = TRUE) |
| Friends-counter | The number of friends an human has | 0 to (number-of-humans—1) |

links become visible when the relation value is higher than the friendship-threshold. All global parameters are described in  Table 2.

## 4    Model Validation and Analysis

Before analysing the model, we validated certain constants to establish reliable settings. These settings were used in the model analysis, which will be discussed further below.

### 4.1    Simulating Real Life Batteries

Two formulas are used to update the social battery: the field-social-battery and the interact-social-battery. Both formulas involve a field-battery-constant and an interact-battery-constant. To calibrate these constants, we performed a sensitivity analysis to find the constants that best represent reality to our experience.

The percentage of humans that do not engage in social interaction is the %COF. We measured the %COF of the introverts and the extraverts and calibrated the field-battery-constant and interact-battery-constant in such a way that with the default parameters, the extraverts mean %COF was around 25% and the mean %COF for the introverts was around 40%.

**Table 2** Overview of global parameters

| Parameter | Use | Value | Default |
|---|---|---|---|
| Status-weight | The importance of status in the status conferral | 0–1 | 0.33 |
| Personality-weight | The importance of personality in the status conferral | 0–1 | 0.33 |
| Relation-weight | The importance of the relationship in the status conferral | 0–1 | 0.33 |
| Volatility* | The importance of the conferral in status changing | 0.1 | 0.1 |
| Sb-threshold | The percentage of social battery needed to interact for introverts | 0–100 | 60 |
| Interact-battery-constant | The rate at which social interaction affects the social battery | 0–30 | 10 |
| Field-battery-constant | The rate at which chilling on the field affects social battery | 0–30 | 30 |
| Friendship-threshold | The value at which people are considered friends | 0–1 | 0.6 |
| Introversion-weight | The importance of introversion to have the energy to socially interact | 0–40 | 20 |
| Number-of-characteristics | The number of characteristics the humans have | 0–20 | 8 |
| Number-of-humans | The number of humans placed on the field | 2–100 | 60 |

*Hard-coded

## 4.2 Resemblance

To determine the number of characteristics to include in the model, we chose the number of characteristics at which the average number of friends stabilised, which was 8.

When the switch adaptive-characteristics? was switched on, within a year almost all characteristics combinations have become the same: the resemblance between the humans is 1 for every relation. Although everyone possessing the same combination of characteristics is not realistic, it is interesting to see that with adaptive characteristics, the humans form more friendships ($p < 0.0001$) (Fig. 2).

## 4.3 Mere Exposure Friendship

One of the interesting things our model showed was that the higher the number of humans in the field, the lower the average number of friendships formed. The highest number of friendships forms at a group size of 20 humans and with more humans, the average number of friends decreases (Fig. 3). This phenomenon can

**Fig. 2** Effect of adaptive characteristics on the average number of friends when the model was run with default settings

be explained by the fact that with more humans, a human has a lower chance of meeting another human they already met before and thus contributing interaction to their relation. Beside this, the decay has a higher effect in the model with more humans, because they meet less repeatedly than in a model with less humans. When two people interact repeatedly, this often results in the two people liking each other, known as the mere-exposure effect [6]. This makes a high number of humans in the model result in relatively few friendships, whereas in a field with a lower number of humans, everyone gets at least one friend. Extraversion, resemblance or status does not matter in that case, the mere-exposure effect increases the relations enough to form friendships. From our own experience, this is also the case in the real world. In a smaller classroom, the chances of becoming friends with the other students are much higher than in a large lecture room.

Interestingly, this effect is larger for introverts than for extraverts. The number of humans at which the extraverts form the most friendships is 20, while for introverts the optimum lies around 10 humans (Fig. 3).

## 4.4 Status Importance

In the formula for calculating the status conferral, the weights of status, personality and previous relationship together is 1. The importance of status can be evaluated by setting the status weight to 1 (making the personality weight and the relation weight 0) and comparing to the results with status weight is 0 (making the personality weight and relation weight both 0.5). Figure 4 shows that the importance of status did not affect the number of friends. Interestingly, the final statuses of the humans did differ between the status weights. When the status was neglected, the final statuses fell in a

**Fig. 3** Effects of group size (number of humans) on the number of friends for introverts, extraverts and the average for all humans when run 5 times



**Fig. 4** The effect of the status weight on the number of friends plotted against the final statuses of the humans

range from 0.2 to 0.7, while with the status being important the final statuses ranged from 0 to 1.

## 5    Discussion and Conclusion

In summary, we created the Friendship Field model to simulate friendship formation including the concept of social battery.

The Friendship Field showed that if humans can adapt their characteristics to become more similar to their friends, the average number of friends increases. Even when the resemblance between all humans is 1, not all humans become friends,

because a relation is not only influenced by resemblance, but also by extraversion, status, appeal and kindness. The adaptivity of the characteristics is a discussion point because in our model the adaptation goes fast compared to our own experience. A future modification to our model could therefore be to include more characteristics and to make people less susceptible to change.

A pattern that resulted from the Friendship Field analysis was the mere-exposure friendship. In a larger group, the probability that one encounters a certain human multiple times is lower than in a smaller group. Therefore, in smaller groups relatively more friendships arise due to the mere-exposure effect. As this is a known phenomenon in social psychology, our model shows an accurate representation of the real world. The optimal group size for friendship formation is for introverts smaller than for extraverts (10 and 20). An explanation for this observation could hold that for introverts, a smaller group size makes it easier to meet people, while the extraverts meet more easily in general so also in larger groups. The optimum for the extraverts is higher because in a larger group, extraverts are able to make more friends due to more possibilities for a friendship. For introverts, having more possibilities for friendships does not benefit them due to their reluctance to meet new people.

The sensitivity analysis suggests that increasing the importance of status in the status conferral results in befriending people with a similar status. When status is neglected, humans can become friends regardless of their status difference. Humans with a low status receive status from other humans with a higher status more easily due to personality and relation, thereby influencing each other's status and forcing the final statuses to be within a smaller range. We expect that when status is important, humans with a high status only become friends with humans with a high status and vice versa, explaining the broad range of final statuses. However, in the current model analysis this effect cannot properly be assessed. Future research might also investigate whether people with a similar status are more likely to become friends. Also, the importance of status and status conferral is dependent on the cultural environment. Our model set the importance of status conferral relatively low because a non-hierarchical environment was assumed. As the importance of status and status conferral is adjustable, the model can serve as a base model for future research in different cultural environments.

Lastly, the friendship threshold in the Friendship Field is the same for every human. However, in real life, the friendship threshold can differ per person. For future research, a different friendship threshold for every human in the model can be incorporated to simulate an even more realistic model.

In conclusion, the Friendship Field proposes a new factor with aspects of Kemper's organism as a reference group theory for friendship formation in social sciences: the social battery. Moreover, it can be used to simulate dyadic friendship formation, based on extraversion, resemblance and status. The model is supported by the theory of the mere-exposure effect on friendship formation. Because of its generality, the Friendship Field can be used for further research in friendship dynamics using the concept of a social battery.

# References

1. Baumeister, R., Leary, M.R.: The need to belong: desire for interpersonal attachments as a fundamental human motivation - PubMed. Psychol. Bull. **117**(3), 497–529 (1995)
2. Bonser, F.G.: Chums: a study in youthful friendships. Pedagogical Seminary **9**(2), 221–236 (1902). https://doi.org/10.1080/08919402.1902.10534181
3. Buhrmester, D.: Intimacy of friendship, interpersonal competence, and adjustment during preadolescence and adolescence. Child Dev. **61**(4), 1101 (1990). https://doi.org/10.2307/1130878
4. Costa, McCrae.: Chapter 19, part 2: the five-factor theory of personality. In: Five Factor Model of Personality (n.d.)
5. Fehr, B.A.: In: Bennett, T.K.: (ed.) Friendship Processes. SAGE Publications, Inc. (1995)
6. Hewstone, M., Stroebe, W., Jonas, K.: An Introduction to Social Psychology, 6th edn. John Wiley & Sons Inc. (2015)
7. Hofstede, G.J., Student, J., Kramer, M.R.: The status–power arena: a comprehensive agent-based model of social status dynamics and gender in groups of children. AI Soc. 1–21 (2018). https://doi.org/10.1007/s00146-017-0793-5
8. Houghton, A.: Understanding personality type: extraversion and introversion. BMJ **329** (2004). https://doi.org/10.1136/sbmj.0411410
9. Kemper, T.D.: Elementary Forms of Social Relations: Status, Power and Reference Groups. Routledge (2017)
10. Marmaros, D., Sacerdote, B.: How do friendships form? Q. J. Econ. **121**(1), 79–119 (2006). https://doi.org/10.1093/QJE/121.1.79
11. Ojanen, T., Sijtsema, J.J., Hawley, P.H., Little, T.D.: Intrinsic and extrinsic motivation in early adolescents' friendship development: friendship selection, influence, and prospective friendship quality. J. Adolesc. **33**(6), 837–851 (2010). https://doi.org/10.1016/J.ADOLESCENCE.2010.08.004
12. Prinstein, M., Dodge, K.: Understanding peer influence in children and adolescents. Choice Rev. Online **46**(03), 46–1775–46–1775 (2008). https://doi.org/10.5860/CHOICE.46-1775
13. PsychReel: What Is Your Social Battery? (A Complete Guide) – PsychReel (2022)
14. Rodkin, P.C., Farmer, T.W., Pearl, R., Acker, R.V.: They're cool: social status and peer group supports for aggressive boys and girls. Soc. Dev. **15**(2), 175–204 (2006). https://doi.org/10.1046/J.1467-9507.2006.00336.X
15. Selfhout, M., Burk, W., Branje, S., Denissen, J., van Aken, M., Meeus, W.: Emerging late adolescent friendship networks and big five personality traits: a social network approach. J. Pers. **78**(2), 509–538 (2010). https://doi.org/10.1111/J.1467-6494.2010.00625.X
16. Urban Dictionary: Urban Dictionary: Social Battery. Urban Dictionary (2017)

# The Importance of Dynamic Networks Within a Model of Politics

**Ruth Meyer** and **Bruce Edmonds**

**Abstract**  Many simulation models of social influence are for the theoretical exploration of the outcomes resulting from certain mechanisms. They therefore tend to be relatively focussed on one mechanism at a time—the KISS approach. Here we take a more KIDS approach, looking at the interaction of two mechanisms within an evidence-led simulation of political behaviour in Austria 2013–2017. In this simulation there is not only the mutual social influence of attitudes (within a 7D space), but this social influence is constrained by a social network. However, one can also allow this social network to adapt based on the interactions between agents, so the social attitudes and social networks co-evolve. In this model, we find that (a) whether the social network is allowed to adapt is more important to the outcomes than the particular kind of social network it is initialized with, but also that (b) (given all the other mechanisms, parameters and structures in this model) a changing social network seems essential to getting outcomes that are qualitatively similar to the patterns in the observed polling data.

**Keywords**  Agent-based simulation · Opinion dynamics · Social influence · Social network

## 1  Introduction

There are now a lot of models that incorporate opinion dynamics, many of them following [1]. For a structured survey of such models, see [2]. Most of these models are intended for the abstract exploration of the consequences of their mechanisms, which is easier if the model is kept relatively simple (certainly free of those details considered not essential to this task). In such models, homophily basically determines who will influence whom—although any two agents can interact, they only influence each other if their opinions are sufficiently close (the difference less than their individual uncertainty). However, when trying to understand social influence

R. Meyer (✉) · B. Edmonds
Centre for Policy Modelling, Manchester Metropolitan University, Manchester, UK
e-mail: meyer.ruth@gmail.com

in observed cases this assumption is not plausible—people only try interacting with a restricted range of people, namely those in their extended social networks (those they interact with face-face, on the phone, on social media etc.). In such cases a social network constrains social influence in addition to homophily as in [3, 4]. This changes the structure of influence, for example a person with opinions that are very different from most others is unlikely to repeatedly contact random others in the hope of finding someone with similar views to their own—more plausibly, they will adapt their social network so that their interactions will be more fruitful more of the time.

In this paper, we look at the importance of the social network on the social influence process within a simulation of political behaviour, specifically some of the politics in Austria between 2013 and 2017. Unlike those models which are aimed at exploring the consequences of abstract mechanisms, this model aimed to be led by the available evidence and data (following the 'KIDS' approach [5]). This results in a much more complicated model. Here, we are interested in the following questions:

1  How does the social network change how social influence works within this model?
2  What social network elements seems to be necessary to get anything like the observed polling outcomes?

## 2 A Model of Voting and Party Competition in Austria

The model used for the exploration of different social influence mechanisms is an agent-based model of voting and party competition in Austria [6]. It simulates the development in Austrian party politics between the national elections of 2013 and 2017, a period that was affected by the refugee crisis of 2015/2016, the ensuing rise of the populist FPÖ, and the leadership-change in and shift to the right by the conservative ÖVP. Parties and voters are agents interacting within a political space spanned by seven policy issues ranging from economical, societal, and environmental topics to immigration policy. Each of these is interpreted as a spatial dimension within a left–right ideological spectrum. Parties and voters take positions on particular issues with lower values indicating they are ideologically left-leaning and higher values indicating they are ideologically right-leaning. The respective values for each voter and party agent are initialised from empirical data: the 2013 Austrian National Election Study (AUTNES) [7] for the voters and the Chapel Hill Expert Survey administered in 2014 [8] for the parties.

Other agent attributes are also defined by the data. For the voters these are demographic characteristics (age, gender, level of education, residential area, income situation), political attitudes (closest party, level of political interest, propensities to vote for any of the parties, probability to vote in the election) and up to three issues of the political space they find most important. For the parties these are their names (included are the seven major Austrian political parties at the time, namely SPÖ,

ÖVP, FPÖ, Grüne, NEOS, BZÖ, and Team Stronach) and equally up to three issues identified as most important. Both parties and voters assign weights to them according to their importance.

Empirical data on issue salience in the public opinion available from the Eurobarometer series of surveys [9] is used as a proxy to model the influence of the media. After matching the relevant Eurobarometer categories to the seven issues represented in the model and rescaling the data, the respective values are applied as probabilities to select the topic to talk about during voter interactions to emulate the media's influence on voter opinion.

The behaviour of voter and party agents is based on theories from the political science literature. To attract voters, parties apply one of a variety of strategies to position themselves in the political landscape [10, 11]; they can choose from "Sticker" (stick to their ideological positions), "Aggregator" (move towards the centre of supporters), "Satisficer" (move like an aggregator until the aspired vote-share is reached), or "Hunter" (seek votes opportunistically by changing direction whenever the vote-share drops). The movement of both "Hunter" and "Aggregator" are restrained to a party's most important issues.

Voters use another set of decision-making strategies to decide which party to vote for. The five strategies identified by [12] comprise "Rational choice", which chooses the party closest in all seven dimensions; "Confirmatory", which picks the party a voter feels closest to (taken from AUTNES); "Fast and frugal", which only looks at the two most important issues to determine the closest party; "Heuristic-based", in which a voter follows recommendations from friends; and "Go-with-your-gut", where voters follow their instinct.

Voters can change their opinions on any of the policy issues due to social influence. This is realized as a bounded confidence opinion dynamics approach, in which randomly paired voter agents only interact if their ideological distance falls under a certain threshold. This threshold represents a voter's 'affective level' and is different for each agent [13]. As the outcome of an interaction voters either move closer together on the discussed topic (agreement) or further apart (disagreement) [14].

The best results obtained with this model using an empirically determined mix of voter decision strategies qualitatively matched the target data, which are the observed opinion polls from 2013 to 2017 (see Fig. 1). However, only a small number of runs did this.

## 3   Social Influence Mechanisms

The social influence mechanism implemented in the Austria model is an opinion dynamics approach which assumes bounded confidence [1, 16]. Two agents only interact if their opinion on the policy issue chosen for discussion (influenced by the Eurobarometer data) is not too dissimilar, i.e., does not exceed a certain confidence threshold. Similarity of opinion is measured as the Euclidean distance between the two agents' opinion in the political space. Following [13], the confidence threshold

**Fig. 1** Observed polling data for the period 2013–2017 [15] versus model-generated polling data

is interpreted as affective involvement and is therefore different for each agent. The Austria model derives this value from empirical data, in particular the political interest of voters, assuming that higher political interest coincides with stronger involvement.

## 3.1 Random Mixing (Totally Connected Network)

In the first version, as usual with opinion dynamics models, interaction partners are chosen randomly from the whole population of agents. In social network terms, this can be interpreted as everyone being connected to everyone else. An analysis of the interactions happening during simulation runs with a small number of discussions per time step (parameter *discussion-freq* set to 1) show that any two randomly chosen

voter agents talk to each other at most 5–6 times over the course of the simulation (208 time steps). However, this is very rare; most will never interact (>70%) or only once (about 22%). Agents have between about 40 to over 400 different interaction partners—numbers at odds with some evidence from political science research, which suggests that the size of political discussion networks is relatively small: people tend to talk to 0–5 other people about politics [17].

## 3.2 Fixed Social Networks

To investigate if the realism of the model can be improved by choosing discussion partners from a voter's social links as suggested by political science research, we consider four different network topologies.

- A regular random network, where each voter is connected to exactly *n* randomly chosen other voters (with *n* specified by model parameter *number-of-friends*).
- An Erdös-Rényi random network, where each configuration of a network with the given mean degree is equally likely; the algorithm used to create this network keeps adding links between randomly chosen pairs of voters until the mean degree (model parameter *number-of-friends*) is reached.
- A scale-free network obtained from preferential attachment, i.e., the probability to connect with a voter rises with the number of links this voter already has.
- The homophily-based network as already implemented in the model, where each voter forms links with other voters most similar in age, education, and residential area from a pool of randomly chosen individuals.

To achieve networks as close as possible to the specification of political discussion networks with 0–5 discussants for every voter, we set the parameter *n* to 3. Table 1 shows the resulting typical values for the different topologies and a population of 1060 voters. The chosen social network is created at model initialisation and remains fixed during a simulation run.

**Table 1** Social network characteristics

| Network type | Total number of links | % voters with 0–5 links | Mean degree | Max. degree | Min. degree | Number of unconnected voters |
|---|---|---|---|---|---|---|
| Homophily-based | 1064 | 99.3 | ≈2 | 7 | 0 | 123 |
| Regular random | 1590 | 100 | 3 | 3 | 3 | 0 |
| Erdös-Rényi | 1590 | 91.3 | 3 | 11 | 0 | 50 |
| Scale-free | 1059 | 95.5 | ≈2 | 58 | 1 | 0 |

### 3.3 Dynamic Social Networks

Keeping the network fixed means that interactions outside the existing social links are not possible. Since these links are still mostly assigned randomly, however, some connections may function less well than others. Some linked voters might be ideologically too far apart on one or more issues for them to ever engage in a conversation on that topic, whereas others might interact but disagree repeatedly. The simulated time frame of four years is also long enough for it to be possible that voters could make new acquaintances to have political discussions with.

We therefore consider an alternative scenario with dynamic networks, where agents may form new random links, friend-of-friend links or drop links with those they disagree with a lot. To this end we introduce three new model parameters: the maximum number of disagreements before the link is dropped (*drop-threshold*), the chance to make a new link (*new-link-prob*) and the proportion for new links to be created with friends of a friend ( *fof-prop*). The outcome of any interaction between two voters is recorded on the link that connects them and stored in a list (-1 for disagreement, $+$ 1 for agreement). At the end of each simulation step, a process to evolve the social network is added. This first deletes all links where the number of disagreements exceeds the drop threshold. Then each voter has the chance to form a new link with either a friend of a friend (80%) or a randomly chosen other voter (20%).

In the experiments reported here, the drop threshold was set to 10 and the probability for a new link to 0.007. While the latter number looks rather small, it avoids an excessive increase in the number of links, keeping the overall 'shape' of the network close to the requirements for political discussion networks.

## 4 Discussion of Results

### 4.1 Effect of Fixed and Dynamic Social Networks

Fixed and dynamic networks are explored through a set number of different scenarios, defined by varying a few chosen model parameters. These govern how often political discussions happen amongst voters (*discussion-freq*: 1, 2 or 5), how easily voters are convinced to change their opinion (*voter-adapt-prob*: 0.5 or 1) and the shape (*network-type*: one of the four different topologies homophily-based, regular random, Erdös-Rényi, preferential attachment) and variability of the social network (model parameter *dynamic-network?* switched off or on). Each scenario is simulated 50 times with the same set of random number seeds.

To compare the different scenarios, we look at election results in the form of possible government formation and measure voter satisfaction as distance to the new government in the two most important issues. Government formation here solely takes the vote shares of parties into account. The largest party forms a coalition with

the next largest party or parties until they reach a majority (> 50%). The ideological positions of such a government in the political space are then computed as the weighted averages of the coalition partners. While this may result in very unrealistic coalitions, for example combining the populist FPÖ with the Greens, it is still a suitable indication of the outcome of a simulation.

To see if voter interaction via the social network improves the realism of the results, each run is compared to the observed historical data. We find that for the fixed networks, none of the runs come close and that the network topology does not make much of a difference. The SPÖ prevails as the biggest party throughout, while the ÖVP comes out as the second biggest party in about 70% of runs, forming a coalition with the SPÖ. The populist FPÖ manages to join the government in up to a third of the cases, mostly in 3-party coalitions. The change of issue salience in the public opinion (rise of the immigration topic) never leads to a dramatic gain for the FPÖ but rather benefits the ÖVP temporarily (see Fig. 2 for an example). This effect is slightly more pronounced with increased discussion frequency, coinciding with a decrease in the government participation of the populists. Figure 3 illustrates the subtle trends with regard to voter interaction.

The results differ for the dynamic networks, i.e., if voter agents are allowed to gain new discussion partners and stop talking to people they disagree with a lot during a simulation. Regardless of network type, there are no longer any 3-party coalitions, and the Greens are never in government. The larger parties win enough vote share to only need one other coalition partner and the Greens are not amongst those. The



**Fig. 2** Typical run with a fixed network (scenario: scale-free network, discussion frequency 2, voter adaptation probability 1)

**Fig. 3** Composition of notional governments over different scenarios across all fixed networks

SPÖ is still either the biggest or the second biggest party, but the FPÖ now manages to win up to 27% of cases depending on the parameter settings defining the voter interaction (*discussion frequency*, *voter-adapt-prob*): starting from 3% (scenarios 1–1, 2–0.5) to 16% (scenario 2–1), to 27% (scenario 5–1). The gain for the ÖVP is even more dramatic, ranging from 3% (scenarios 1–1, 2–0.5) to 43.5% (scenario 5–1). The more people talk and convince each other, the higher the chance that the FPÖ or ÖVP become the largest party instead of the SPÖ (see Fig. 4).

Change in issue salience in the public opinion now has a noticeable effect, though the advantage is still mostly taken by the ÖVP. A few runs do come qualitatively close to the historical data and here the network type does make a difference: while



**Fig. 4** Composition of notional governments over different scenarios across dynamic networks

**Fig. 5** Best model results with dynamic networks

the Erdös-Rényi and Regular Random network both display examples of "successful" runs the other two network types (homophily-based and scale-free) do not. Figure 5 shows the best result, obtained with the Erdös-Renyi network in scenario 5–1 (*discussion-freq* 5, *voter-adapt-prob* 1).

## 4.2 Sensitivity to Network Type and Dynamics

To see the impact of various non-network settings in fixed and dynamic network cases we varied the following parameters with 10 independent runs for each set (192,000 simulation runs in total).

- *discussion-freq*: {1, 2}
- *max-p-move*: {0.5, 1}
- *voter-adapt-prob*: {0.5, 1}
- *max-salience-change*: {1.5, 3}
- *dynamic-network?*: {true, false}
- *network-type*: {"homophily-based political discussion network", "regular random network", "preferential attachment", "Erdös-Rényi random network"}
- *number-of-friends*: {1, 3, 5}

   The key output measure of interest we use is the level of voter satisfaction with the notional elected government at the end of a simulation run, i.e., the proportion of

**Fig. 6** Overall contrast of dynamic versus non-dynamic networks for each of four initial network configurations

voters within 10% of the centroid of government policies in their two most important issues. In each such diagram the error bars show one standard deviation either way (Fig. 6).

The significance of the dynamism of the network is evident. For all four different topologies used to initialize the interaction network it clearly makes a difference whether the network evolves during a simulation or not.

The sub-case where it made the most difference was with the following settings (Fig. 7):

- *max-salience-change*: 3
- *voter-adapt-prob*: 1
- *max-p-move*: 1
- *discussion-freq*: 2

The sub-case where it made the least difference was as follows (Fig. 8):

- *max-salience-change*: 1.5
- *voter-adapt-prob*: 0.5
- *max-p-move*: 0.5
- *discussion-freq*: 1

Lower values of the model parameters *max-salience-change*, *voter-adapt-prob*, *max-p-move* and *discussion-freq* result in less difference between dynamic and non-dynamic networks, but this is still a clearly identifiable difference.

**Fig. 7** The sub-case where the dynamism of the network made the greatest difference



**Fig. 8** The sub-case where the dynamism of the network made the least difference

# 5  Conclusion

Our main conclusion is that, given all the other features of the model (many but not all of which were suggested by the available evidence), the network dynamics are essential for producing results like that of the reference case in this model. The opinion dynamics and network dynamics co-evolve and reinforce each other. This echoes the results in some previous models with social influence aimed at understanding political processes, namely the abstract model described in [18] and a more evidence-led model looking at the reasons why people bother to vote [19]. The dynamism of the network makes most difference with more discussion between agents, more adaptivity by voters in terms of attitudes and salience change, and a greater adaptivity from those parties who change policies in response to the voter attitude landscape. Since many models are for theoretical exploration, they tend to focus on either social influence of attitudes or adaptation of social networks. This work suggests that, at least in some empirically-driven cases, both mechanisms might be needed, since they can act to amplify or dampen each other's effects.

# References

1. Deffuant, G., Neau, D., Amblard, F., Weisbuch, G.: Mixing beliefs among interacting agents. Adv. Comp. Syst. **3**(01n04): 87–98. (2000)
2. Flache, A., Mäs, M., Feliciani, T., Chattoe-Brown, E., Deffuant, G., Huet, S., Lorenz, J.: Models of social influence: towards the next frontiers. J. Artific. Soc. Soc. Simul. **20**(4), 2. (2017). http://jasss.soc.surrey.ac.uk/20/4/2.html, https://doi.org/10.18564/jasss.3521
3. Das, A., Gollapudi, S., Munagala, K.: Modeling opinion dynamics in social networks. In: Proceedings of the 7th ACM International Conference on Web Search and Data Mining (WSDM '14), pp. 403–412. Association for Computing Machinery, New York, NY, USA (2014). https://doi.org/10.1145/2556195.2559896
4. Stauffer, D., Sousa, A., Schulze, C.: Discretized opinion dynamics of the deffuant model on scale-free networks. J. Artific. Soc. Soc. Simul. **7**(3) 7 (2004). https://www.jasss.org/7/3/7.html
5. Edmonds, B., Moss, S.: From KISS to KIDS – an 'anti-simplistic' modelling approach. In: Davidsson, P., et al. (eds.) Multi Agent Based Simulation 2004, vol. 3415, pp. 130–144. Springer, Lecture Notes in Artificial Intelligence (2005)
6. Meyer, R., Fölsch, M., Dolezal, M., Heinisch, R.: An evidence-driven model of voting and party competition in Austria. In: Czupryna, M., Kamiński, B. (eds.) Advances in Social Simulation. Proceedings of the 16th Social Simulation Conference, 20–24 Sept 2021, pp. 261–273 (2022)
7. Kritzinger, S., Zeglovits, E., Aichholzer, J., Glantschnigg, C., Glinitzer, K., Johann, D., Thomas, K., Wagner, M. (2017): AUTNES Pre- and Post Panel Study 2013. GESIS Data Archive, Cologne. ZA5859 Data file Version 2.0.1. https://doi.org/10.4232/1.12724
8. Polk, J., Rovny, J., Bakker, R., Edwards, E., Hooghe, L., Jolly, S., Koedam, J., Kostelka, F., Marks, G., Schumacher, G., Steenbergen, M., Vachudova, M., Zilovic, M.: Explaining the salience of anti-elitism and reducing political corruption for political parties in Europe with the 2014 Chapel Hill expert survey data. Res. Polit. **4**(1), 1–9 (2017)
9. The European Commission's Eurobarometer Surveys. Available online at https://www.gesis.org/en/eurobarometer-data-service/home
10. Laver, M.: Policy and the dynamics of political competition. Am. Polit. Sci. Rev. **99**(2), 263–281 (2005)

11. Muis, J., Scholte, M.: How to find the 'winning formula'? Conducting simulation experiments to grasp the tactical moves and fortunes of populist radical right parties. Acta Politica **48**(1), 22–46 (2013)
12. Lau, R., Kleinberg, M., Ditonto, T.: Measuring voter decision strategies in political behavior and public opinion research. Public Opin. Q. **82**(S1), 911–936 (2018)
13. Schweighofer, S., Garcia, D., Schweitzer, F.: An agent-based model of multi-dimensional opinion dynamics and opinion alignment. Chaos: Interdisc. J. Nonlinear Sci. **30**(9), 093139 (2020)
14. Baldassarri, D., Bearman, P.: Dynamics of political polarization. Am. Soc. Rev. **72**(5), 784–811 (2007)
15. Opinion Polls in Austria from 2013 to 2017. https://en.wikipedia.org/wiki/Opinion_polling_for_the_2017_Austrian_legislative_election
16. Hegselmann, R., Krause, U.: Opinion dynamics and bounded confidence: models, analysis and simulation. J. Artifi. Soc. Soc. Simul. **5**(3), 2 (2002). https://www.jasss.org/5/3/2.html
17. Lake, R., Huckfeldt, R.: Social capital, social networks, and political participation. Polit. Psychol. **19**(3), 567–584 (1998)
18. Edmonds, B.: Co-developing beliefs and social influence networks – towards understanding socio-cognitive processes like Brexit. Qual. Quant. **54**(2), 491–515 (2020). https://doi.org/10.1007/s11135-019-00891-9
19. Lafuerza, L.F., Dyson, L., Edmonds, B., McKane, A.J.: Staged models for interdisciplinary research. PLoS ONE **11**(6), e0157261 (2016). https://doi.org/10.1371/journal.pone.0157261

# Management and Economics

# A Preliminary Study of Individual Based Crowd Simulation Based on Bayesian Nash Equilibrium

**Yiyu Wang** , **Jiaqi Ge** , **and Alexis Comber**

**Abstract** The lack of experimental datasets for individual behaviours has hindered the systematic studies of pedestrian behaviours as well as the refined development of regular laws of individual movement in simulation models. This research developed a simulation model for crowd evacuation on the basis of Bayesian Nash Equilibrium (BNE) and a Multi-Agent System (MAS). BNE was introduced in this research to augment the rationality of individual decision-making process in evacuation simulation and to assist pedestrians in discovering an optimal evacuation route to avoid congestions. A series of simulation experiments were conducted to evaluate the performance of the initial model, and the current experimental results demonstrate a noticeable positive influence of BNE on reducing evacuation time. A detailed introduction of the establishment and implementation details of the model as well as model analysis have been provided in this paper. Limitations and a few optional research directions in the future are also discussed.

**Keywords** Agent based modelling · Bayesian nash equilibrium · Emergency evacuation · Pedestrian behaviours · Multi-agent system.

## 1 Background

The rising prevalence of large crow gatherings in public places has attracted the attentions of relevant researcher who have sought to discover efficient measures for crowd evacuation to minimise the risk of causalities in unexpected incidents such as disasters [1]. However, a series of obstacles, such as the lack of empirical experimental dataset, have hindered further studies on pedestrian flow to some extent. To overcome these challenges, a range of field observations and theoretical studies have

Y. Wang (✉) · J. Ge · A. Comber
School of Geography, University of Leeds, Leeds LS2 9JT, UK
e-mail: gyywa@leeds.ac.uk

been carried out to gain insights into the movements of pedestrian flow during emergency evacuations in difference scenarios [2, 3]. Despite these efforts, the systematic features of pedestrian movements remain poorly understood, highlighting the demand for realistic experimental data on individual behaviours during emergency evacuations in different scenarios.

To address this gap, this study proposes an effective simulation model for pedestrian flow in which individuals enable to evacuate more realistically and will be provided suggestions for faster evacuation routes. To achieve this, Bayesian Nash Equilibrium (BNE) was employed within a Multi-Agent System (MAS) to augment the simulation of individual decision-making process during emergency evacuations. The proposed model seeks to improve the realism of pedestrian behaviours during emergency evacuations in order to address the shortage of realistic experimental datasets for relevant research on pedestrian flow.

The model development process was outlined using the ODD + D protocol proposed by Müller et al. [4], which emphasizes the importance of specifying the objectives, design concepts, and details of the simulation model. The analysis of simulation results has also been provided in the following sections to evaluate the performance of the proposed model.

## 2   Design Concept

### 2.1   Theoretical Background

Bayesian Nash Equilibrium (BNE) is a widely adopted concept in game theory, which extends the standard framework of Nash Equilibrium by introducing the possibility of incomplete information [5]. In some real-world scenarios, players may not have access to complete information about the strategies and payoffs of other players in the game. Instead, they need to choose a strategy to play only depending on their beliefs about the others' strategies based on their experience and knowledge of the game. BNE offers a more comprehensive framework for analysing the interactions of strategies taken by different players in such scenarios, by considering the uncertainty that arises from incomplete information. Consequently, the capability of BNE to incorporate the impact of incomplete information and the changes of belief status over time is particularly important in the analysis of game-theoretic scenarios involving complex and dynamic environments.

BNE provides a way to update the probabilities of the other players' strategies on the basis of new evidence. And players using BNE rules are capable to maximize their own expected utility by updating their beliefs about the strategies played by their neighbours in light of the new information obtained [5]. That is, the crucial factor in determining the next strategy to play is the probability distribution of other players making the same decision in the same game. In this model, the rules of BNE

is reflected in the calculations of relevant utilities during agents' decision-making process, which will be explained in details in the following section.

## 2.2 Individual Decision-Making

The process of how the agent competes with its neighbours and select the moving direction is described in detail below.

BNE provides a relative flexible approach to modelling the decision-making of strategies, by allowing for the consideration of the multiple factors which may affect the decision-making of players in the game. In this model, there are no restrictions on the sequence of an agent's decision-making, offering a more realistic description of the dynamic and complex environments involving the interactions between players as well as strategies [6].

In this model, *Total Utility* ($U_t$) is associated with the distance from the current location to the exit (i.e. *Distance Utility* $U_d$), the number of nearby agents in same patch (i.e. *Comfort Utility* $U_c$) and the probabilities of neighbours moving to the same patch in the next time step (i.e. *Expected Comfort Utility* $U_{ec}$). The sum of $U_d$ and $U_{ec}$ is used to calculate the value of total utility, as expressed in Eq. 1.

Reverse movement has not been allowed in the current implementation of the model which means that agents are restricted to choose from six available directions $P_0$, $P_1$, …, $P_5$ at evert time step (refer to Fig. 1). It is highly probable that multiple agents may converge in the same patch during evacuation, resulting in congestions with varying degrees. To address this issue, BNE was introduced here to enable agents to avoid congestions and evacuate efficiently, which can be achieved by calculating the probabilities of the number of agents entering each patch in the next time step. -Each BNE agent will assess the total utilities of six optional patches in front of itself (see Fig. 1) and selects the one with the maximum total utility for movement in the next time step.

$$U_t = U_d + U_{ec} \tag{1}$$

Agents in this model use the information on neighbours' current location and their probability distribution of moving directions to predict future situation of crowd gatherings. And in this study, all the related BNE utilities were set as patch attributes. *Total utility* used for simulation is related to three key elements: *Distance Utility*, *Comfort Utility* and *Expected Comfort Utility*, and the detailed identifications of these utilities are provided in the followings.

**Distance Utility**. It reflects the distance from current position to the exit and keeps increasing as the exit is closer, as shown in Eq. (2).

$$U_d = (D - d)/D \tag{2}$$

**Fig. 1** Optional directions of agents in the model

where, $d$ represents the distance to the exit and $D$ refers to the length of the diagonal of simulation space.

**Comfort Utility**. It is a crucial component of the model which makes it possible to capture the comfort levels experienced by the agents situated in the same patch. A series of coefficients has been involved to provide a numerical reflection of individual comfort in the given space. Considering the real-world scenarios where limited space capacity may negatively affect the comfort level of residents, comfort utility in this model would be assigned a value of 0 when the number of agents occupying the same patch exceeds four. This, to some extent, ensures that the model is able to capture the realistic constraints of physical space and provide a support for the following evacuation simulations.

$$U_c = \begin{cases} 1.00, n \le 2 \\ 0.51, n = 3 \\ 0.07, n = 4 \\ 0.00, n \ge 5 \end{cases} \tag{3}$$

where, n refers to the number of agents in one patch.

**Expected Comfort Utility**. It is the multiplication of comfort utility $U_c$ and the probability of agents in the patch in the next time step $p(n)$. Specially, all agents on the patch where they are on and their 8 neighbour patches have been taken into account of the calculation of expected comfort utility.

$$U_{ec} = \sum_{n=0}^{4} C_N^n P_m^n (1 - P_m)^{N-n} \tag{4}$$

where, $n$ refers to the number of agents in one patch; $N$ represents the total number of agents on 10 patches; $P_m$ means the probability of the agent moving to the patch, which can be adjusted by the corresponding slider in the model.

The relationship of these BNE utilities is illustrated in Fig. 2.

**Fig. 2** The schema of BNE utilities

# 3 Implementation Details and Model Analysis

## 3.1 Implementation Details

The initial model of this study was implemented in Java using NetLogo. The source code has been published on COMSES and available at https://doi.org/10.25937/75wf-aa82 [7].

By default, model world was initialized with 2000 agents scattered randomly over the evacuation space and presumed that individuals could evacuate through the two exits on either side with a width of 6. The initial speed of individuals has been determined based on the number of their neighbours in 8 surrounding patches. The simulation environment representing the physical space was divided into 1360 (68*20) patches, and each patch allowed over one agent to occupy. For each patch, both distance utility and expected comfort utility were calculated at the initialization stage. Agents following BNE behavioural model tended to choose a patch where was in front of them and had the highest value of total utility to move towards. The agents' decisions, as well as their speed and other parameters, was continuously updated per time step until the end of simulation.

In addition, the model allows for the regulation of several relevant variables such as the percentage of agents using BNE, exit size, and so on, through the corresponding sliders. The moving pattern of agents can be selected using the chooser "moving-pattern" which enables the observation of the behaviours and movements of evacuating pedestrians in different scenarios.

Four options for moving patterns are provided in this model, which are Random Follow (RF), Shortest Route (SR), BNE mixed with RF and BNE mixed with SR. To be specific, the Shortest Route option means that all the agents move directly to the exit which may result in congestions and extend the evacuation time; Random Follow allows agents randomly select a neighbour in their views to follow and update their

following targets every time step; The BNE mixed with Shortest Route and BNE mixed with Random Follow options allow for a specific percentage of agents using BNE behavioural model, which can be adjusted by the slider "Percentage-of-agents-with-BNE", to participate in evacuation simulations, while the rest follow either the SR or RF patterns.

It is important to note that no input data is currently used in this model and all the parameters are set by default or regulated by the corresponding sliders.

## 3.2   Preliminary Results

NetLogo BehaviorSpace was introduced in this study to appraise whether and how pedestrian evacuation could be impacted by BNE involved. A series of experiments were carried out with three different behavioural models: BNE, Random Follow, Shortest Route). For each pattern, 50 runs were undertaken with 100% agents evacuated successfully in each simulation. The evacuation time for agents with different moving patterns was also recorded in order to evaluate the performance of the BNE model during emergency evacuations.

Figure 3 presents a comparison of the evacuation time among three different moving patterns: Shortest Routes, Random Follow, and BNE at 100% adoption rate. The results demonstrate that the evacuation time when all agents use BNE to evacuate is apparently lower than that of other two patterns, which means that the adoption of BNE function has a significant impact on reducing exit time during simulations.

The findings of this study point out the potential of using BNE in realistically simulating evacuating behaviours and improving pedestrian evacuation strategies in emergency scenarios, and incorporating BNE as a decision-making model could effectively reduce evacuation time and improve the overall efficiency of evacuation efforts, which can ultimately contribute to more comprehensive and safety measures for crowd management.

## 4   Conclusion and Future Work

This study aims to develop an evacuation simulation model for pedestrian flow, which incorporates Bayesian Nash Equilibrium (BNE) to enhance the rationality of pedestrian decision-making within a Multi-Agent System (MAS) that simulates pedestrian movements and behaviours during evacuating. A series of experiments have been conducted under different moving patterns in order to evaluate the performance of the proposed model. The current experimental results demonstrate a clear positive influence of BNE on reducing evacuation time, suggesting the effectiveness of the proposed model.

**Fig. 3** Evacuation time of three moving patterns: 100%BNE, random follow, and shortest route (number of agents: 2000)

However, it should be noticed that the proposed model, while offering valuable insights into individual decision-making in emergency evacuation scenarios, has still been subject to several limitations that need to be addressed to improve its applicability and predictive accuracy.

Firstly, certain model attributes, such as moving speed, comfort utility, and so on, require further calibration in order to accurately reflect the individual movements and behaviours in the real-world evacuation scenarios. And sensitivity analysis is also required to assess the influences of these attributes on the performance of the model.

Secondly, as several variables, such as the number of pedestrians, the percentage of BNE users, and so on, were held fixed during simulations in this paper, additional simulation experiments still need to be conducted with a wider range of parameter configurations, including different widths of exits, combinations of different moving patterns, a broader range of the number of agents, to gain a relatively comprehensive assessment of the effects of BNE on evacuation process.

Thirdly, real or pre-existing evacuation datasets that could provide valuable empirical evidence need to be employed in order to validate the reliability and predictive accuracy of the proposed model.

The role of space played in individual decision-making process need to be further investigated through the introduction of different types of blockades into the simulation space before or during evacuations, which would also enable researchers to examine how pedestrians adapt their decisions in response to unexpected barriers

or changes in the simulation environment as well as evaluate the influences of such factors on the overall evacuation process.

Furthermore, the incorporation of greater self-organizing behaviours among pedestrians, such as competitive behaviours, could make the model more adaptable to a broader range of relevant research on pedestrian simulations, and also contribute to better simulating individual behaviours and movements in complex and dynamic environments.

In conclusion, the proposed model provides valuable insights into individual decision making in emergency evacuation scenarios and offers a promising avenue for improving crowd evacuation strategies, yet several limitations need to be addressed through further research. The integration of BNE within a MAS presents a novel approach to simulating individual decision-making process and improving the realism and accuracy of evacuation simulations of pedestrian flow, which could ultimately contribute to more effective implementations for crowd management.

# References

1. Babojelić, K., Novacko, L.: Modelling of driver and pedestrian behaviour–a historical review. Promet-Traffic Transp. **32**(5), 727–745 (2020)
2. Feng, Y., Duives, D., Daamen, W., Hoogendoorn, S.: Data collection methods for studying pedestrian behaviour: a systematic review. Build. Environ. **187**, 107329 (2021)
3. Rozo, K.R., Arellana, J., Santander-Mercado, A., Jubiz-Diaz, M.: Modelling building emergency evacuation plans considering the dynamic behaviour of pedestrians using agent-based simulation. Safety Sci. **113**, 276–284 (2019)
4. Müller, B., Bohn, F., Dreßler, G., Groeneveld, J., Klassert, C., Martin, R., Schlüter, M., Schulze, J., Weise, H., Schwarz, N.: Describing human decisions in agent-based models–ODD+ D, an extension of the ODD protocol. Environ. Model. Softw. **48**, 37–48 (2013)
5. Ui, T.: Bayesian Nash equilibrium and variational inequalities. J. Math. Econ. **63**, 139–146 (2016)
6. Wang, Y., Ge, J., Comber, A.: An agent-based simulation model of pedestrian evacuation based on bayesian nash equilibrium. J. Artif. Soc. Soc. Simul. **26**(3), 6 (2023). https://doi.org/10.18564/jasss.5037
7. Wang, Y., Ge, J., Comber, A.: An agent-based simulation model of pedestrian evacuation based on Bayesian Nash Equilibrium (Version 1.0.0). CoMSES Computational Model Library. https://doi.org/10.25937/75wf-aa82. Last accessed 06 Jul 2022
8. Wang, Y., Ge, J., Comber, A.: An evacuation simulation model of pedestrian flow using Bayesian Nash equilibrium and a multi-agent system. AGILE: GISci. Series **3**, 1–5 (2022)

# Agent-Based Simulations and Process Mining: A Green BPM Case Study

**Emilio Sulis**

**Abstract**  This paper investigates the interactions between agent-based modeling and process mining, which is an increasingly widespread applied discipline in the context of business process management. In particular, we explore a practical "green BPM" perspective. We propose a simulation approach to describe the environmental effects of two different health policy strategies: the hospital scenario and the home hospitalisation scenario. Furthermore, we demonstrate the application of process discovery techniques from an 'event log' obtained from an agent-based simulation. As a case study, we propose the home hospitalisation service of a hospital in Turin. The traditional hospital-centred care scenario shows how much pollution is produced, compared to the scenario with the same number of patients treated at home. Finally, we describe how discovering processes from event logs opens the way to improving the organisation's service management.

**Keywords**  Agent-based modeling · Process mining · Event-logs

## 1  Introduction

The discipline of agent-based computational management has gained interest in both research and industrial fields [1, 2]. The growth of process-oriented and data-driven applications makes it possible to achieve different research perspectives through the integration of techniques and tools available in the various disciplines. Recently, an important research area focused on business processes management (BPM) has extensively explored the adoption of event-logs (from real or simulation data) by means of pattern extraction, as well as techniques for comparing the process model with log data, in the process mining (PM) approach [3]. Applications can be multiple, such as management tools, monitoring, decision support in different sectors, e.g. in industry, environment, or health.

E. Sulis (✉)
Computer Science Department, University of Turin, Turin, Italy
e-mail: emilio.sulis@unito.it

Research in economics and business has typically focused on discrete-event simulations or system dynamics [4]. The most recent approach involves agent-based modelling (ABM), which facilitates the exploration of the complexity of emerging phenomena and stakeholder understanding. In addition to simulation, the exploration of event-logs makes it possible to examine processes of interest (e.g. sequences of industrial activities, environmental events, patient care pathways in healthcare) by means of discovery algorithms. This paper explores the integration of agent-based simulations and PM, exploring a case of sustainability in health policy, with a practical application of *green BPM*. The aim is to compare the impact of pollution generated by a health care management based on the home hospitalisation model instead of the traditional health care model with centralisation of patients in hospital. Using pollutant emissions as an indicator to compare environmental impact, our result shows the benefits of the innovative home hospitalisation service (HHS) compared to the impact of hospital-only care of the same patients. In addition, the agent-based simulation of HHS can be used to generate an *event-log* of the care processes. The log of activities can be explored with process discovery techniques to demonstrate the usefulness of the PM perspective.

In the remainder of the article, Sect. 2 describes the background with related work, the methodology, and the case study. Simulation results are presented in Sect. 3, while Sect. 4 describes a process discovery application. Section 5 concludes the article by introducing some future work.

## 2  Background

*Related work.* Agent-based simulation has been applied to healthcare in several directions, e.g. in a management perspective [5], epidemiology for influence-like virus spreading [6], social problems as drug addiction or elderly facilities [7]; spatial impact (spatial perspective, closed environment) [8].

In the business process life cycle, business process modelling is a central phase of analysis. The BPM research area typically explores modeling processes by means of standard notations such as Petri Nets, BPMN, or Direct Follower Graph (DFG) [9]. Nevertheless, agent-oriented approaches recently concerned operational business processes from real-world event data [10], as well as agent-based business process simulation [11].

The adoption of real event-logs for process modeling in PM implies two main applications, i.e. the discovery of the workflow model from real data or *process discovery* [12], as well as the comparison between the ideal model and the discovered one or *conformance checking* [13]. event-logging systems in healthcare are still not widespread, requiring costly investments in sensors and system architectures [14]. In cases where no real logs are available, an alternative is to run the simulations from which the event-log is derived [15]. Our case study investigates an healthcare application in a perspective of green BPM [16], by focusing on the environmental degradation caused by pollutant emissions from vehicles.

*Hospital-at-home service case study.* An application of interest concerns, in particular, the HHS of the main hospital in Turin, the City of Health and Science. This service has been running for over 35 years, at the Complex University Geriatrics Unit of the Molinette Hospital. A medical-nursing team takes care of patients at home on a daily basis, with benefits for both the family and healthcare.[1] The HHS is a care model that has proved to be a valid alternative to hospitalisation for a whole range of acute and chronic exacerbated diseases, with a lower risk of complications such as infections, delirium and malnutrition. During the COVID-19 pandemic emergency there was a need to reduce the pressure on hospital beds and a consequent reshaping of the HHS service for the management of mild COVID-19 patients but above all for the management of non-infected patients, also supported by telemedicine tools [17].

*Methodology.* The main parameters of the agent-based simulation has been addressed by a secondary analysis of the real data for the last 5 years to reconstruct the most frequent patterns, identify the arrival frequencies, the patients' needs. In addition, the model was built by interacting with domain experts, i.e. doctors from the HHS. We analysed the data on patients' examinations for each of the work teams involving nurses and doctors. We obtained the frequency of each procedure, integrated with qualitative concerns by domain experts. E.g. 'midline' medication is performed in about the 40% of patients, then typically repeated at least 3–4 times. Finally, we consider the order of the sequence of visits that each patient must always follows, e.g. structured visit, then medical visit, then nurse visit.

To compare the simulation results, the *HHS scenario* considers the 5 teams (vehicles) moving to reach about 700 patients in one year, with about 20–24 procedures each day. The same initial setting addresses a simulation of the traditional care scenario, where patients are treated in the hospital (HOS scenario). In this case, a number of relatives/caregiver/parents move to the hospital during the patient's stay. We considered a relative for each patient moving to the hospital (by own vehicle) each day. This causes a negative impact on environmental conditions. As a key performance indicator, we adopt a measure of 30 grammes of $CO_2$ per kilometre (g $CO_2$/km) to compute the pollution of vehicles.

To explore ABM, we rely on NetLogo, a platform suitable to address a small-scale simulation, as in our case. We adopt two extensions: *GIS* to import the shapefile corresponding to the streets of the city of Turin, on which the vehicles move in the simulations, and *time*, to obtain realistic data for constructing the event-log in the simulation results. We adopt a subset of real patient addresses from the HHS service in order to obtain a more realistic dataset on which to base the simulation. The code is available on the repository associated to the paper.[2] Figure 1 describes the agents for teams and patients in a GIS environment.

Finally, the PM exploration requires the generation and extraction of an event-log from the simulation that must include traces of all activities of each patient from the initial commitment to the HHS ('TAKE-CHARGE') to the last activity with the discharge ('DISMISSION'). Each trace in the log has three main relevant three

---

[1] HHS details at: https://www.ospedaleadomicilio.it/.

[2] https://github.com/sulemil/SSC22.

**Fig. 1** The simulation of patients and teams moving in a GIS-based representation

attributes: the case identifier of the patient (anonymously labelled by an ID code), the event (name of the procedure), and the timestamp (at the precision of minutes). The simulation can also record the completion date, but for a demonstration exercise, as in this case, the start date is sufficient. In addition, we also record the team that performed the activity (added as an attribute of the event-log). The complete event-log obtained during the simulation has been exported in a CSV file. A validation of the log with domain experts is necessary, to focus on more interesting cases (e.g. we filter out cases with only one activity and lasting less than three days). On the top of the event-log, we apply process discovery techniques to investigate the sequence of activities of patients. We adopt PM Fluxicon's Disco tool,[3] by using the fuzzy miner algorithm [18].

## 3 Scenario Analysis

*Hospital-at-home service scenario.* The first scenario of interest is that of home hospitalisation (HHS scenario). Results for pollution refer to a period of one year. The execution of 100 simulations gives an average value of total pollution of 5,457,427.8 g $CO_2$/km (with a standard deviation of 119,365.6). Figure 2 describes the model with an example of the generation (on the left) of the event log for later exploration with process discovery techniques (as detailed in Sect. 4).

*Hospital scenario.* In the case of the traditional scenario with hospital care (HOS scenario), the simulation shows an increased presence of pollution mainly due to the movements of relatives and caregivers to reach the hospital. Despite having set a relatively low number of visits (one per day for the duration of hospitalisation) the emission values reach very high values. In fact, the average total pollution for 100 executions is 52,211,149.8 g $CO_2$/km (standard deviation of 7,239,572.1). Figure 2 describes on the right part of the interface with the HOS scenario, while the central area describes a visualization of the pollution as clouds.

---

[3] https://fluxicon.com/disco/.

**Fig. 2** A view of the simulation model for hospital-at-home service (HHS) and hospital (HOS) scenario. In the left area, the construction of the event-log. On the right, some indicators about pollution and healthcare management.

## 4 Discovery Results

### 4.1 Event-Log

The simulation results can be used to generate a synthetic dataset, i.e. an event log including the case ID of the patient, the operation, and the starting date.

```
45,TAKE-CHARGE,2022-01-19 13:23,Team33
67,MEDICATION-MIDLINE,2022-01-19 13:35,Team30
45,TRASF-1,2022-01-19 13:41,Team31
67,MIDLINE,2022-01-20 10:11,Team33
18,MEDICAL-VISIT,2022-01-20 10:21,Team32
18,ENEMA,2022-01-20 10:22,Team30
45,ECOGRAPHY,2022-01-19 13:11,Team31
```

In the example, three patients (numbers 45, 67 and 18) were involved in different medical operations at the corresponding start date, with the hospital team that performed it. The complete event-log can be inspected to remove possible outliers, and filter appropriately before proceeding. We also consider a warm-up period of one month. For example, we removed incomplete patient traces (e.g. those that do not end at the time of the end of the simulation). Then, the event log is analysed with a

PM tool to identify the main features in an initial exploration. The traces in the log corresponds to 20 activities for 479 cases, with a median duration of 10 days. These value has been validated by domain experts as well as by the historical data of HSS.

## 4.2 Control-Flow and Performance Analysis

Based on the event-log described above, discovery techniques can obtain the control-flow perspective of the activities of patients. Figure 3 describes the output from DISCO of the sequence of activities (the number indicates the frequency of each activity in the log). The visualization (a DFG diagram) allows to easily interpret data. Moreover, modifications can be made on the simulated model, once validated by domain experts, to check whether the activity flow of the log reflects, improves or worsens the prospects of interest.

A further analysis makes it possible to check times in the process flow, highlighting the arcs that take the longest, on average, between one activity and another. In this way, possible bottlenecks can emerge, to be studied with the domain experts in order to propose effective solutions, perhaps even through the automation of certain tasks or punctual scheduling. Figure 4 describes the output of process discovery from the event-log including only the 50% of activities, to focus on more frequent activities.



**Fig. 3** The control flow perspective describes the sequence of activities and their frequency, as registered in the event-log

**Fig. 4** The performance perspective shows the average duration between the corresponding tasks, clearly indicated by the weighted arcs. This is useful to identify bottlenecks in the flow of activity.

The weight of the arcs indicates in most cases a time between 5 and 7 days between the tasks, while some activities are closer in time. These tools can be used, by changing the settings appropriately, for management and monitoring purposes.

## 5   Conclusions

We explored the possible synergy between an agent-based simulation and process mining techniques, by exploring a practical application in green BPM. Limitations are in the inherent complexity of the task, while the use case presented here with Net-Logo is illustrative and demonstrative of the framework to be used. The introductory examples in the modeling and simulation platform open the way to further investigation for performing analysis with different techniques or directions, e.g., parameter sweeping, learning algorithms, social network analysis, conformance checking. Finally, an interesting possibility lies in the automatic construction of agent-based simulations from real event-logs.

# References

1. Wall, F.: Agent-based modeling in managerial science: an illustrative survey and study. Rev. Manag. Sci. **10**(1), 135–193 (2016)
2. Wall, F., Leitner, S.: Agent-based computational economics in management accounting research: opportunities and difficulties. J. Manag. Account. Res. **33**(3), 189–212 (2021)
3. van der Aalst, W.M.P.: Process Mining—Data Science in Action, 2nd edn. Springer, Berlin (2016)
4. Djanatliev, A.: Hybrid simulation for prospective healthcare decision-support: system dynamics, discrete-event and agent-based simulation. Ph.D. thesis, University of Erlangen-Nuremberg (2015)
5. Cabrera, E., Taboada, M., Iglesias, M.L., Epelde, F., Luque, E.: Optimization of healthcare emergency departments by agent-based simulation. In: Sato, M., Matsuoka, S., Sloot, P.M.A., van Albada, G.D., Dongarra, J.J. (eds.) Proceedings of ICCS. Procedia Computer Science, vol. 4, pp. 1880–1889. Elsevier (2011)
6. Marini, M., Brunner, C., Chokani, N., Abhari, R.S.: Enhancing response preparedness to influenza epidemics: agent-based study of 2050 influenza season in Switzerland. Simul. Model. Pract. Theory **103**, 102091 (2020)
7. Keyes, K.M., Shev, A., Tracy, M., Cerdá, M.: Assessing the impact of alcohol taxation on rates of violent victimization in a large urban area: an agent-based modeling approach. Addiction **114**(2), 236–247 (2019)
8. Rajabi, M., Pilesjö, P., Shirzadi, M.R., Fadaei, R., Mansourian, A.: A spatially explicit agent-based modeling approach for the spread of cutaneous leishmaniasis disease in central Iran, Isfahan. Environ. Model. Softw. **82**, 330–346 (2016)
9. Dumas, M., La Rosa, M., Mendling, J., Reijers, H.: Fundamentals of Business Process Management, vol. 1, 2nd edn. Springer, Berlin (2018)
10. Tour, A., Polyvyanyy, A., Kalenkova, A.A.: Agent system mining: vision, benefits, and challenges. IEEE Access **9**, 99480–99494 (2021)
11. Sulis, E., Taveter, K.: Agent-Based Business Process Simulation—A Primer with Applications and Examples. Springer, Berlin (2022)
12. Augusto, A., Conforti, R., Dumas, M., Rosa, M.L., Maggi, F.M., Marrella, A., Mecella, M., Soo, A.: Automated discovery of process models from event logs: review and benchmark. IEEE Trans. Knowl. Data Eng. **31**(4), 686–705 (2019)
13. Jans, M., Weerdt, J.D., Depaire, B., Dumas, M., Janssenswillen, G.: Conformance checking in process mining. Inf. Syst. **102**, 101851 (2021)
14. Sulis, E., Amantea, I.A., Aldinucci, M., Boella, G., Marinello, R., Grosso, M., Platter, P., Ambrosini, S.: An ambient assisted living architecture for hospital at home coupled with a process-oriented perspective. J. Ambient Intell. Humanized Comput. 1–19 (2022)
15. van der Aalst, W.M.P.: Process discovery from event data: relating models and logs through abstractions. WIREs Data Min. Knowl. Discov. **8**(3) (2018)
16. Couckuyt, D., Looy, A.V.: A systematic review of green business process management. Bus. Process. Manag. J. **26**(2), 421–446 (2020)
17. Amantea, I.A., Arnone, M., Leva, A.D., Sulis, E., Bianca, D., Brunetti, E., Marinello, R.: Modeling and simulation of the hospital-at-home service admission process. In: Obaidat, M.S., Ören, T.I., Szczerbicka, H. (eds.) Proceedings of SIMULTECH, pp. 293–300. SciTePress (2019)
18. Günther, C.W., Rozinat, A.: Disco: discover your processes. In: Lohmann, N., Moser, S. (eds.) Proceedings of the Demonstration Track of the 10th International Conference on Business Process Management (BPM 2012), Tallinn, Estonia, 4 Sept 2012. CEUR Workshop Proceedings, vol. 940, pp. 40–44. CEUR-WS.org (2012)

# Collaborative Search and Autonomous Task Allocation in Organizations of Learning Agents

**Stephan Leitner**

**Abstract** This paper introduces a model of multi-unit organizations with either static structures, i.e., they are designed top-down following classical approaches to organizational design, or dynamic structures, i.e., the structures emerge over time from micro-level decisions. In the latter case, the units are capable of learning about the technical interdependencies of the task they face, and they use their knowledge by adapting the task allocation from time to time. In both static and dynamic organizations, searching for actions to increase the performance can either be carried out individually or collaboratively. The results indicate that (i) collaborative search processes can help overcome the adverse effects of inefficient task allocations as long as there is an internal fit with other organizational design elements, and (ii) for dynamic organizations, the emergent task allocation does not necessarily mirror the technical interdependencies of the task the organizations face, even though the same (or even higher) performances are achieved.

**Keywords** *NK* framework · Adjacent walk · Evolutionary organizational design · Guided self-organization

## 1 Introduction

Designing organizations includes a multiplicity of decisions, such as breaking down the task of the larger problem for smaller units, allocating responsibility and authority to departments and individuals, coordinating behavior through incentives, communication, leadership, and routines, among others, and it is well known that an organization's design substantially impacts the organization's performance [3, 4]. The main challenges of organizational design are to achieve an external fit, i.e., to design organizations for dynamic and uncertain situations and perhaps even situations that have not been seen before [3], and an internal fit among the organizational design elements

S. Leitner (✉)
University of Klagenfurt, Klagenfurt, Austria
e-mail: stephan.leitner@aau.at

[13], which might be particularly difficult when organizations evolve through phases of their life-cycle and the employees' capabilities and knowledge are dynamic [5].

There are two main world-views on organizational design: First, classic approaches follow the premise of the rational actor and postulate that organizational design is the result of *deliberate* decisions [14]; following this view, managers design feasible organizations top-down. Second, evolutionary approaches consider that organizational structures emerge bottom-up. The latter approach includes a shift from the macro-level to the micro-structures, focusing on mechanisms that drive the emergence of organizational design elements [6]. This paper addresses two such micro-level issues: First, limited information, learning, and adaptation, and second, collaborative search processes.

*Limited information, learning, and adaptation* concern the technical characteristics and decomposition of the task the organization faces. Previous research recommends that an organization's structure should mirror the task's technical interdependencies (mirroring hypothesis) [12]. There are ambiguous results regarding this hypothesis; some previous research criticizes it based on empirical evidence, and, at the same time, there are also empirical results that support it [1, 10]. Efficiently designing organizations top-down and in line with the mirroring hypothesis requires that the technical structure (i.e., the structure of interdependencies) is public knowledge. In reality, this structure is unknown and unclear in most cases [11]. Highly complex tasks might not only be challenging to decompose; previous research argues that increasing the number of interdependencies also unfolds non-linear effects that lead to performance drops, what is often labelled as 'complexity catastrophe' [7]. This paper addresses both cases of organizational design mentioned before; there are scenarios in which (i) the technical interdependencies of the task are known beforehand, and organizations are designed top-down, and (ii) the technical interdependencies are not known, but agents learn about it over time and can adapt the task allocation over time.

This paper relies on situated learning theory to model *collaborative search processes*, according to which search processes might take place in interactive communities [16]. While traditional search algorithms mainly focus on individual search processes [15], this paper enriches the models of an organization with distributed and autonomous decision-makers by a social network that constitutes organizational connections. These connections are then used to autonomously coordinate search behavior, resulting in collaborative search efforts. For dynamic and static organizations, the paper tests whether there are organizational design elements, such as control mechanisms and (collaborative) search processes, that either reinforce or weaken the 'complexity catastrophe'.

The remainder of this paper is organized as follows: Sect. 2 introduces the model and the method of data analysis, Sect. 3 presents and discusses the results. Finally, Sect. 4 summarizes and concludes the paper.

## 2  Model

The model builds on the well-known *NK* framework [15]. The organization comprises $M \in \mathbb{N}$ organizational units, referred to as agents henceforth. All agents face an $N$-dimensional decision problem with $K$ interdependencies among them, where $N \in \mathbb{N}$ and $K \in \mathbb{N}_0$. The interdependencies shape the decision problem's complexity. Due to limited capacities, the agents cannot solve the entire decision problem alone, but they decompose it into $M$ sub-problems that agents can handle (Sect. 2.1). The agents aim to increase their utilities by employing an individual or collaborative search processes (Sect. 2.2). The agents know that they face a complex decision problem. However, they do not know the actual number and structure of interdependencies between decisions. Still, they are endowed with the capability to learn about the structure of interdependencies (Sect. 2.3). Also, the agents use their knowledge by adapting the task allocation from time to time (Sect. 2.4). For $t = \{1, \ldots, T\} \subset \mathbb{N}$ periods it is observed how the agents' decisions affect the organization's performance. The model was implemented in Matlab® (R2022a).

### 2.1  Task Environment and Decomposition

The decision problem faced by the agents consists of $N$ binary decisions and is formalized by $\mathbf{d} = [d_1, d_2, \ldots, d_N]$, where $d_n \in \{0, 1\}$ and $n = \{1, \ldots, N\} \subset \mathbb{N}$. Every decision $d_n$ contributes $f(d_n) \sim U(0, 1)$ to the organization's performance. Due to interdependencies among decision, the performance contribution $f(d_n)$ might not only be affected by decision $d_n$ but also by $K$ other decisions. The corresponding contribution function for decision $d_n$ is formalized by $f(d_n) = f(d_n, d_{i_1}, \ldots, d_{i_K})$, where $\{i_1, \ldots, i_K\} \subseteq \{1, \ldots, n-1, n+1, \ldots, N\}$ and $0 \leq K \leq N - 1$. The organizations' performance is the average of all performance contributions:

$$P(\mathbf{d}) = \frac{1}{|\mathbf{d}|} \sum_{n=1}^{|\mathbf{d}|} f(d_n). \tag{1}$$

The agents are limited in their capabilities and/or resources, i.e., they might have limited cognitive capacities, limited time, or limited further resources to solve the decision problem. Consequently, they have to collaborate to find a feasible solution to the complex decision problem captured by the task environment. To do so, they decompose the decision problem into $M$ sub-problems $\mathbf{d_m}$, where $m = \{1, \ldots, M\} \subset \mathbb{N}$ and $[\mathbf{d}_1, \ldots, \mathbf{d}_M] = \mathbf{d}$. For agent $m$, the decisions $\mathbf{d}_m$ represent the area of responsibility, while the complement $\mathbf{d}_{-m} = \mathbf{d} \setminus \mathbf{d}_m$ is referred to as residual decisions. The agents can observe the solutions to their sub-problem $\mathbf{d}_m$ at any time. However, the solutions to the residual decision problem $\mathbf{d}_{-m}$, can only be observed *after* implementation.

**Decisions** (K=2, Decomposable structure)

| Performance contributions | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | x | x | x | | | | | | | | | | | | |
| 2 | x | x | x | | | | | | | | | | | | |
| 3 | x | x | x | | | | | | | | | | | | |
| 4 | | | | x | x | x | | | | | | | | | |
| 5 | | | | x | x | x | | | | | | | | | |
| 6 | | | | x | x | x | | | | | | | | | |
| 7 | | | | | | | x | x | x | | | | | | |
| 8 | | | | | | | x | x | x | | | | | | |
| 9 | | | | | | | x | x | x | | | | | | |
| 10 | | | | | | | | | | x | x | x | | | |
| 11 | | | | | | | | | | x | x | x | | | |
| 12 | | | | | | | | | | x | x | x | | | |
| 13 | | | | | | | | | | | | | x | x | x |
| 14 | | | | | | | | | | | | | x | x | x |
| 15 | | | | | | | | | | | | | x | x | x |

**Decisions** (K=5, Non-decomposable structure)

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | x | x | x | | | | x | | | x | | | x | | |
| 2 | x | x | x | | | x | | | x | | | x | | | |
| 3 | x | x | x | | | | | x | | | x | | | x | |
| 4 | | | | x | x | x | | | | x | | | x | | x |
| 5 | x | | | x | x | x | | | | | x | | | x | |
| 6 | | x | | x | x | x | | | | | | x | | | x |
| 7 | | x | | | | | x | x | x | | | | x | | x |
| 8 | x | | | | x | | x | x | x | | | x | | | |
| 9 | | x | | | x | | x | x | x | | | | | x | |
| 10 | | x | | | | x | | | | x | x | x | | | x |
| 11 | x | | | | x | | | x | | x | x | x | | | |
| 12 | | x | | | | x | | | x | x | x | x | | | |
| 13 | | x | | | | x | | | x | | | | x | x | x |
| 14 | | | x | | | x | | | x | | | | x | x | x |
| 15 | | | x | | | x | | | x | | | | x | x | x |

**Fig. 1** Interdependence matrices

This paper considers two stylized interdependence structures presented in Fig. 1, where an 'x' indicates that a decision and a performance contribution are interdependent. The task allocation indicated by black lines is used for scenarios with *top-down designed organizations*. The considered structures are of complexity $K = 2$ and $K = 5$, representing a fully decomposable and non-decomposable decision problem, respectively.

In *organizations with emergent structures*, the allocation of tasks to agents might be adapted from time to time, i.e., agents might swap tasks. In these scenarios, in every period $t \bmod \tau = 0$, agents can adjust the task allocation (Sect. 2.4). In contrast, in periods $t \bmod \tau \neq 0$, agents seek to maximize their utility given the currently active task allocation (Sect. 2.2), where $\tau \in \mathbb{N}$. The task allocation in period $t = 1$ follows a random process that allocates tasks equally so that the number of decisions assigned to agent $m$ is $|\mathbf{d}_m| = N/M$.

## 2.2 Utility Functions and Search Processes

The performance contributions of agent $m$'s own ($\mathbf{d}_{mt}$) and residual decisions ($\mathbf{d}_{-mt}$) in $t$ are denoted by $P(\mathbf{d}_{mt})$ and $P(\mathbf{d}_{-mt})$, respectively. The organization employs a linear outcome-based incentive scheme that shapes the agents' utility functions. In particular, the parameter $\alpha \in \mathbb{R}^+$ is used to weight the agents' own and residual performances, respectively, where $0 \leq \alpha \leq 1$. Agent $m$'s utility at period $t$ is formalized by

$$U(\mathbf{d}_{mt}, \mathbf{d}_{-mt}) = \alpha \cdot P(\mathbf{d}_{mt}) + (1 - \alpha) \cdot P(\mathbf{d}_{-mt}). \tag{2}$$

The agents seek to maximize their utilities by employing one of the following two variants of a hill-climbing algorithm:

**Individual search**. Agent $m$ discovers a solution $\mathbf{d}^*_{mt}$ to their partial decision problem in period $t$ characterized by a Hamming distance of 1 to the solution $\mathbf{d}_{mt-1}$, i.e., $\mathbf{d}^*_{mt}$ is different from $\mathbf{d}_{mt-1}$ in exactly one position. Direct communication between agents is omitted in individual hill-climbing, so agent $m$ has no information about the other agents' decisions but relies on the other agents' decisions from the previous period, $\mathbf{d}_{-mt-1}$, to compute the utility. Agent $m$ selects the solution to be implemented in $t$ from their options $\mathbf{D}^{\text{ind}}_t = \{\mathbf{d}_{mt-1}, \mathbf{d}^*_{mt}\}$ according to the following rule:

$$\mathbf{d}_{mt} = \arg\max_{\mathbf{d}' \in \mathbf{D}^{\text{ind}}_t} U\left(\mathbf{d}', \mathbf{d}_{-mt-1}\right). \tag{3}$$

**Collaborative search**. Agents are connected in a ring network, and they interact with one of their nearest neighbors with probability $\mathbb{P}$. If they interact, agents $m$ and $n$ jointly perform adjacent hill-climbing [16] to maximize their *joint utility*. They share information about the solutions $\mathbf{d}_m$ and $\mathbf{d}_n$ to their partial decision problem. Let us denote the solutions to the decisions outside the two agents' areas of responsibility by $\mathbf{d}_{-(m,n)} = \mathbf{d} \setminus (\mathbf{d}_m \cup \mathbf{d}_n)$. Then, the agents' joint utility in period $t$ is the mean of the individual utilities in Eq. 2:

$$U^{\text{adj}}\left(\mathbf{d}_{mt}, \mathbf{d}_{nt}, \mathbf{d}_{-(m,n)t}\right) = \frac{1}{2} \cdot (U(\mathbf{d}_{mt}, \underbrace{\mathbf{d}_{-mt}}_{\mathbf{d}_{-(m,n)t} \cup \mathbf{d}_{nt}}) + U(\mathbf{d}_{nt}, \underbrace{\mathbf{d}_{-nt}}_{\mathbf{d}_{-(m,n)t} \cup \mathbf{d}_{mt}})) \tag{4}$$

The two agents discover and share with their counterparts the new solutions $\mathbf{d}^*_{mt}$ and $\mathbf{d}^*_{nt}$. Again, the newly discovered solutions are characterized by a Hamming distance of 1 to the corresponding solutions in the previous period. For the decisions outside their areas of responsibility, the agents $m$ and $n$ rely on the residual solutions implemented in the last period, $\mathbf{d}_{-(m,n)t-1}$. The agents jointly choose the solutions to be implemented in period $t$ from the tuples $\mathbf{D}^{\text{adj}}_t = \{(\mathbf{d}_{mt-1}, \mathbf{d}_{nt-1}), (\mathbf{d}^*_{mt}, \mathbf{d}_{nt-1}), (\mathbf{d}_{mt-1}, \mathbf{d}^*_{nt})\}$ according to the rule

$$(\mathbf{d}_{mt}, \mathbf{d}_{nt}) = \arg\max_{(\mathbf{d}'_m, \mathbf{d}'_n) \in \mathbf{D}^{\text{adj}}_t} U^{\text{adj}}\left(\mathbf{d}'_m, \mathbf{d}'_n, \mathbf{d}_{-(m,n)t-1}\right). \tag{5}$$

**Computation of the overall solution**. The solution to the decision problem that is implemented in period $t$ is the concatenation of the decisions made by all $M$ agents, $\mathbf{d}_t = [\mathbf{d}_{1t}, \ldots, \mathbf{d}_{Mt}]$, and the performance achieved by the organization in $t$ is $P(\mathbf{d}_t)$ (Eq. 1).

## 2.3   Learning Mechanism

The agents know that they face a complex decision problem, but they do not know the
exact structure of interdependencies among decisions. However, agents are endowed
with beliefs on the interdependencies, and they update them in all periods $t \bmod \tau \neq$
$0$. We formalize agent $m$'s belief on the interdependencies between decisions $i$ and $j$
in period $t$ by $b_{mt}^{ij} \in \mathbb{R}$, where $i, j = \{1, \ldots, N\} \subset \mathbb{N}$, $i \neq j$, and $0 \leq b_{mt}^{ij} \leq 1$. The
beliefs $b_{mt}^{ij}$ are computed as the mean of the Beta distribution $B(p_{mt}^{ij}, q_{mt}^{ij})$. For the
initial beliefs, $p_{m1}^{ij} = q_{m1}^{ij} = 1$ so that $b_{m1}^{ij} = 0.5$. During the observation period, agent
$m$ makes decisions in their area of responsibility and fixes the decisions $\mathbf{d}_{mt}$ to be
implemented in $t$ by either following the individual (Eq. 3) or adjacent hill-climbing
algorithm (Eq. 5). If agent $m$ decides to change a decision so that $\mathbf{d}_{mt} := \mathbf{d}_{mt}^{*}$, the
beliefs on interdependencies are updated as follows:

1. Let us denote the decision that has been flipped by agent $m$ in $t$ by $i$, where $d_{it} \in$
   $\mathbf{d}_{mt}$. After implementing the decisions $\mathbf{d}_{mt}$, agent $m$ observes the performance
   contributions of all decisions within their area of responsibility.
2. Whenever agent $m$ observes that the performance contribution of decision $j$
   changes from period $t - 1$ to period $t$ if the decision $i$ is flipped, $p_{mt}^{ij}$ is increased
   by 1, otherwise $q_{mt}^{ij}$ is increased by 1:

$$
\left( p_{mt}^{ij}, q_{mt}^{ij} \right) =
\begin{cases}
\left( p_{mt-1}^{ij} + 1, q_{mt-1}^{ij} \right) & \text{if } f(d_{jt}) \neq f(d_{jt-1}), \\
\left( p_{mt-1}^{ij}, q_{mt-1}^{ij} + 1 \right) & \text{otherwise} .
\end{cases}
\tag{6}
$$

3. Agent $m$ recomputes the beliefs $b_{mt}^{ij}$.

Please note that agents can only observe the performance contributions *within their*
areas of responsibility. Suppose the decision problem is decomposed so that there
are interdependencies with decisions from *outside* an agent's area of responsibility;
in that case, there might be external influence on performance contributions that the
agent cannot identify as such.

## 2.4   Task Re-allocation Mechanism

In all periods $t \bmod \tau = 0$, agents are granted the possibility to re-organize the task
allocation.[1] To account for limitations in resources, every agent is characterized by a
maximum capacity $C_m$ that indicates the maximum number of decisions that agent $m$
can handle at a time. $C_m$ can be interpreted in terms of maximum cognitive capacity
or maximum financial resources, time, manpower, etc., that are available to solve
decision problems.

---

[1] Please note that a re-allocation of decision also affects the computation of the agent's utility in
terms of what is regarded as own and residual performance (see Eq. 2).

**Computation and exchange of signals**. Agents follow the idea of the mirroring hypothesis and aim at maximizing the interdependencies within their own areas of responsibility. The process is organized as follows:

1. Agent $m$ identifies the task $i$ in their own area of responsibility that is associated with the minimum belief on internal interdependencies. Agent $m$ also sends a signal (Eq. 7) that is used as a threshold for trading this decision, i.e., the task is only re-allocated if the other agents' signals exceed the threshold signal.

$$\rho_{mt}^i = \min_{\forall i: d_{it} \in \mathbf{d}_{mt}} \left( \frac{1}{|\mathbf{d}_{mt}| - 1} \sum_{\substack{\forall j: d_{jt} \in \mathbf{d_{mt}} \\ j \neq i}} b_{mt}^{ij} \right) \tag{7}$$

2. Agent $m$ informs the other agents that the task $i$ that fulfils Eq. 7 and the threshold signal $\rho_{mt}^i$. Agents $r$ proceed with the next step and send signals iff $|\mathbf{d}_{rt}| < C_r$.
3. Agents $r$ submit the average belief on the interdependencies between the offered task $i$ with the decisions within his or her area of responsibility $\mathbf{d}_{rt}$ as a signal in period $t$. Agent $r$'s signal for decision $i$ in $t$ is formalized by

$$\bar{\rho}_{rt}^i = \frac{1}{|\mathbf{d}_{rt}|} \sum_{\forall j: d_{jt} \in \mathbf{d_{rt}}} b_{rt}^{ij} \tag{8}$$

**Task re-allocation**. Once all agents sent their signals, for every offer $i$, there are at most $M - 1$ signals. Recall, agent $m$ offered task $i$ at a threshold signal of $\rho_{mt}^i$ and the other agents sent signals $\bar{\rho}_{rt}^i$. Let us denote the set of signals for task $i$ in period $t$ by $\mathbf{P}_t^i$, the maximum signal for task $i$ in period $t$ by $\bar{\rho}_{r*t}^i = \max_{\bar{\rho}_{rt}^i \in \mathbf{P}_t^i}(\bar{\rho}_{rt}^i)$, and the agent sending the maximum signal by $r^*$. The tasks are (re-)allocated as follows: If the maximum signal $\bar{\rho}_{r*t}^i$ is equal to or exceeds the threshold signal $\rho_{mt}^i$, the task $i$ is re-allocated from agent $m$ to agent $r^*$ according to

$$\mathbf{d}_{mt} := \mathbf{d}_{mt-1} \setminus \{d_{it-1}\} \text{ and} \tag{9a}$$
$$\mathbf{d}_{r*t} := \left[ \mathbf{d}_{r*t-1}, d_{it-1} \right], \tag{9b}$$

where $\setminus$ indicates the complement. If the maximum signal $\bar{\rho}_{r*t}^i$ does not exceed the threshold $\rho_{mt}^i$, agent $m$ remains responsible for task $i$, so that $\mathbf{d}_{mt} := \mathbf{d}_{mt-1}$.

## 2.5 Parameters and Data Analysis

**Parameters**. The main parameters are summarized in Table 1. This paper puts particular emphasis on the analysis of the relation between task performance (as the dependent variable) and task complexity $K$, collaborative search probability $\mathbb{P}$, and the incentive parameter $\alpha$ (the independent variables). To assure comparability across

**Table 1** Parameters

| Type | Variables | Notation | Values |
|---|---|---|---|
| Independent variables | Task complexity | $K$ | {3, 5} |
| | Time steps | $t$ | {1 : 1 : 150} |
| | Collaborative search probability | $\mathbb{P}$ | {0 : 0.05 : 0.5} |
| | Incentive parameter | $\alpha$ | {0, 25, 0.5, 0.75} |
| Dependent variable | Normalized task performance | $\tilde{P}(\mathbf{d_t})$ | [0, 1] |
| Other parameters | Number of decisions | $N$ | 15 |
| | Agents | $M$ | 5 |
| | Agents' cognitive capacities | $C_m$ | 5 |
| | Task re-allocation | $\tau$ | {∅, 25} |
| | Number of simulations | $S$ | 800 |

simulation runs, the observed performance $P(\mathbf{d}_{ts})$ is normalized by the maximum attainable performance in that scenario, $P(\mathbf{d}_s^*)$, so that $\tilde{P}(\mathbf{d}_{ts}) = P(\mathbf{d}_{ts})/P(\mathbf{d}_s^*)$. In addition to cases in which agents can adapt the task allocation in every $\tau = 25$ periods, i.e., *emergent organizational structures*, there are benchmark scenarios in which the initial allocation of tasks already follows the mirroring hypothesis (which is indicated the bold lines in Fig. 1) and the agents cannot re-allocate tasks ($\tau = \emptyset$), i.e., *top-down designed organizations*.

**Regressions and partial dependencies**. To analyze the functional dependencies between the dependent and the independent variables, regression neural networks are trained, and partial dependencies are computed [2, 9]. Let $\mathbf{X}$ be the set of all independent variables included in Table 1. The subset $\mathbf{X}^s$ includes the independent variable(s) that are in the scope of the analysis, and $\mathbf{X}^c$ consists of the complementary set of $\mathbf{X}^s$ in $\mathbf{X}$. Then, $f(\mathbf{X}) = f(\mathbf{X}^s, \mathbf{X}^c)$ represents the trained regression model. The partial dependence of the performance on the independent variables in scope is defined by the expectation of the performance concerning the complementary independent variables so that

$$f^s(\mathbf{X}^s) = E_c(f(\mathbf{X}^s, \mathbf{X}^c)) \approx \frac{1}{V} \sum_{i=1}^{V} f(\mathbf{X}^s, \mathbf{X}_{(i)}^c), \tag{10}$$

where $V$ is the number of independent variables in $\mathbf{X}^c$ and $\mathbf{X}_{(i)}^c$ is the $i$th element. By marginalising over the independent variables in $\mathbf{X}^c$, we get a function that depends only on the independent variables in $\mathbf{X}^s$.

**Task allocation efficiency**. The efficiency of task re-allocation is evaluated using the following metric: Let $C(\mathbf{d}_{mt})$ be a count-function that returns the number of interdependencies *within* agent $m$'s sub-problem in $t$. Then, the following ratio of interdependencies within agent $m$'s sub-problem (nominator) to the total number of times the decisions assigned to agent $m$ affect performance contributions (denominator) is used as the task re-allocation efficiency metric:

$$\eta_{mt} = \frac{C(\mathbf{d}_{mt})}{|\mathbf{d}_{mt}| \cdot K} \tag{11}$$

## 3  Results

### 3.1  Effects of Complexity, Time, and Collaborative Search on Performance

**Complexity**. The partial dependencies of performance on complexity are plotted in Fig. 2. The results indicate that whether or not endowing the agents with the capability to re-allocate tasks reinforces the 'complexity catastrophe' [7] depends on the incentive system effective in the organization. In particular, the results for top-down designed organizations reflect the finding that higher levels of complexity result in lower task performance [8]. The results for emergent organizational structures show that individualistic incentives reinforce the effect of complexity on performance. In contrast, task re-allocation appears to slightly weaken (or, at least, not reinforce) this effect in cases with altruistic incentives. Thus, focusing on complexity only, bottom-up designed organizations are best off if they employ altruistic incentives, whereas individualistic incentives result in the most significant drop in performance.

**Time and collaborative search probability**. The partial dependencies of performance on time and collaborative search probability are presented in Fig. 3; top-down and bottom-up organizational designs are indicated by solid and dashed lines, respec-
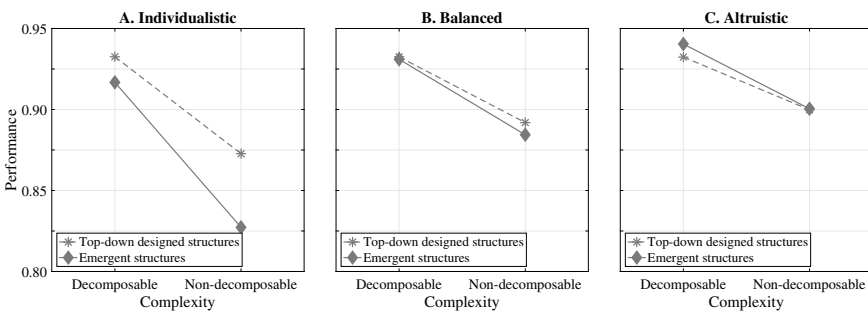


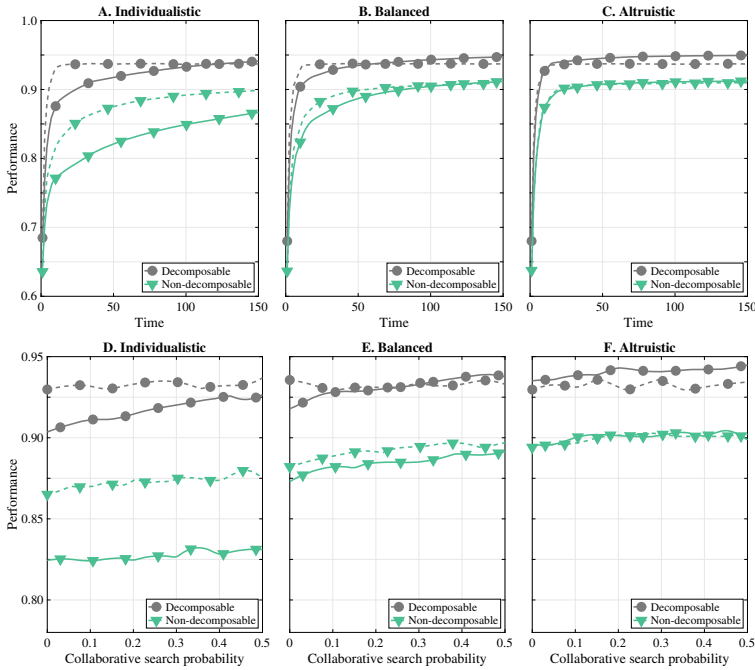**Fig. 2** Partial dependence of performance on complexity

tively. Dots circles ( ● ) indicate scenarios with decomposable tasks, and triangles (▼) stand for non-decomposable tasks.

For *decomposable tasks*, the partial dependencies indicate that the performances in top-down designed organizations grow relatively fastly and reach the upper limit early in the observation period. For emergent organizational structures, both the speed and the upper limit are affected by the incentive parameter: The performance grows relatively slowly and eventually reaches the upper limit of the performance in top-down organizations if *individualistic incentive systems* are effective (Fig. 3a). The partial dependencies of task performance on the collaborative search probability (Fig. 3d) indicate that this pattern is reinforced if the collaborative search probability is low, i.e., the distance between the performances in the two cases gets larger. In contrast, the performances become more similar if the collaborative search probability is high. In the case of *balanced incentive systems* (Fig. 3b), the dependence of the performance on time is relatively similar to panel A, and the collaborative search probability appears not to significantly affect the slopes of the performance curves (Fig. 3e). If *altruistic incentive systems* are effective (Fig. 3c), the performance reacts more substantially to time when the organizational structure is dynamic. The performance is eventually higher compared to the performance in top-down designed organizations. The results presented in Fig. 3f indicate that this effect is reinforced if the collaborative search probability increases. This means that relatively high collaborative search probabilities pay off in performance if altruistic incentive schemes are effective in the organization.

The patterns observed for scenarios with *non-decomposable tasks* are similar to those for decomposable tasks, whereby, as already evident from Fig. 2, relatively lower performances are achieved. For *individualistic incentive schemes* (Fig. 3a), the performance increases faster and reaches a higher level in top-down designed organizations than in cases with emergent structures; increasing the collaborative search probability in these cases only has negligible effects. The performance increases faster but has approximately the same upper limit if *balanced and altruistic incentive mechanisms* are effective in the organization (Fig. 3b, c). When altruistic incentive systems are effective, the performances in top-down and bottom-up designed organizations become very similar; the partial dependencies plotted in Fig. 3f indicate that this pattern is robust against variations in the collaborative search probability.

## 3.2   Task Allocation Efficiency

This section analyses to what extent the emerging organizational structure in scenarios with task re-allocation conforms to the task allocation suggested by the mirroring hypothesis (i.e., the solid lines in Fig. 1). The following task allocation efficiency in scenarios with top-down structures are used as a benchmark: In the case of decomposable decision problems, all interdependencies are internalized into the agents' decision problems (Fig. 1, $K = 2$), and, in consequence, the benchmark efficiency metric reaches a value of 1. For non-decomposable decision problems, only a sub-

Dashed $(--)$ and solid lines $(—)$ stand for benchmark scenarios and scenarios with task re-allocation, respectively.

**Fig. 3** Partial dependence of performance on time and collaborative search probability



**Fig. 4** Cumulative distributions of the task allocation efficiency metric

set of the interdependencies can be internalized; only 6 out of 15 interdependencies (40%) are inside an agents' decision problems in Fig. 1, $K = 5$, and, in consequence, the benchmark efficiency metric is 0.4.

The cumulative distributions of the task allocation efficiency metric are plotted in Fig. 4 (for all agents and all periods). Interestingly, in only approx. 10% of the cases, agents achieve a task allocation efficiency of 0.5 out of 1 in the case of decomposable tasks and 0.3 out of 0.4 for non-decomposable tasks. Even though the signals for task re-allocation are based on the agents' beliefs on interdependencies, the incentive parameter affects the task allocation efficiency: Irrespective of task complexity, altruistic incentive schemes result in a slightly higher task allocation efficiency; this might be driven by an indirect effect coming from the individual search behavior induced by altruistic incentives as well as the resulting update of beliefs on interdependencies.

## 4 Conclusions

This paper presents a model of either dynamic or static organizations, in which search processes are carried out individually or collaboratively. The results indicate that collaborative search processes can indeed weaken the adverse effects of emergent task allocations that do not conform to the mirroring hypothesis. However, this is only true if there is a fit between the search processes and the remaining organizational design elements, namely with the inventive scheme: The results indicate that emergent approaches to organizational design work best with rather altruistic incentive schemes. Surprisingly, the results also indicate that organizations are better off if they follow an emergent design approach together with altruistic incentives if tasks are decomposable: In these cases, the performance even exceeds that of top-down organizations. Thus, the results indicate that the long standing finding that an organization's structure should mirror the technical interdependencies of the task the organization faces is not necessarily applicable in organizations with emergent structures.

This work can be seen as the first step toward an organizational design theory in dynamic organizations with autonomous agents. Further research could, for example, analyze different strategies for task re-allocation (e.g., different ways to compute the signals), different network structures for organizational links, and the effects of collaborative search in networks of organizations. Also, future research might take into account other forms of performance landscapes (e.g., plateaued landscapes).

## References

1. Baldwin, C., MacCormack, A., Rusnak, J.: Hidden structure: using network methods to map system architecture. Res. Policy **43**(8), 1381–1397 (2014)
2. Blanco-Fernandez, D., Leitner, S., Rausch, A.: Dynamic groups in complex task environments: to change or not to change a winning team? (2022). arXiv preprint arXiv:2203.09157
3. Burton, R.M., Obel, B.: The science of organizational design: fit between structure and coordination. J. Organ. Des. **7**(1), 1–13 (2018)

4. Burton, R.M., Obel, B., Håkonsson, D.D.: Organizational Design. A Step-by-Step Approach, 4th ed. Cambridge University Press (2020)
5. Cardinal, L.B., Sitkin, S.B., Long, C.P.: Balancing and rebalancing in the creation and evolution of organizational control. Organ. Sci. **15**(4), 411–431 (2004)
6. Joseph, J., Baumann, O., Burton, R., Srikanth, K.: Reviewing, revisiting, and renewing the foundations of organization design. In: Organization Design. Emerald Publishing Limited (2018)
7. Kauffman, S.A., et al.: The Origins of Order: Self-organization and Selection in Evolution. Oxford University Press, USA (1993)
8. Leitner, S., Wall, F.: Multiobjective decision making policies and coordination mechanisms in hierarchical organizations: results of an agent-based simulation. Sci. World J. (2014)
9. Patel, M.H., Abbasi, M.A., Saeed, M., Alam, S.J.: A scheme to analyze agent-based social simulations using exploratory data mining techniques. Complex Adapt. Syst. Model. **6**(1), 1–17 (2018)
10. Querbes, A., Frenken, K.: Grounding the "mirroring hypothesis": towards a general theory of organization design in new product development. J. Eng. Technol. Manage. **47**, 81–95 (2018)
11. Raveendran, M., Silvestri, L., Gulati, R.: The role of interdependence in the micro-foundations of organization design: task, goal, and knowledge interdependence. Acad. Manage. Ann. **14**(2), 828–868 (2020)
12. Sanchez, R., Mahoney, J.T.: Modularity, flexibility, and knowledge management in product and organization design. Strateg. Manage. J. **17**(S2), 63–76 (1996)
13. Thompson, J.D., Zald, M.N., Scott, W.R.: Organizations in Action: Social Science Bases of Administrative Theory. Routledge (2017)
14. Tsoukas, H.: Organizations as soap bubbles: an evolutionary perspective on organization design. Syst. Prac. **6**(5), 501–515 (1993)
15. Wall, F., Leitner, S.: Agent-based computational economics in management accounting research: opportunities and difficulties. J. Manage. Acc. Res. **33**(3), 189–212 (2021)
16. Yuan, Y., McKelvey, B.: Situated learning theory: adding rate and complexity effects via Kauffman's NK model. Nonlinear Dyn. Psychol. Life Sci. **8**(1), 65–101 (2004)

# Controlling Replication via the Belief System in Multi-unit Organizations

**Ravshanbek Khodzhimatov** [ID]**, Stephan Leitner** [ID]**, and Friederike Wall** [ID]

**Abstract** Multi-unit organizations such as retail chains are interested in the diffusion of best practices throughout all divisions. However, the strict guidelines or incentive schemes may not always be effective in promoting the replication of a practice. In this paper we analyze how the individual belief systems, namely the desire of individuals to conform, may be used to spread knowledge between departments. We develop an agent-based simulation of an organization with different network structures between divisions through which the knowledge is shared, and observe the resulting synchrony. We find that the effect of network structures on the diffusion of knowledge depends on the interdependencies between divisions, and that peer-to-peer exchange of information is more effective in reaching synchrony than unilateral sharing of knowledge from one division. Moreover, we find that centralized network structures lead to lower performance in organizations.

**Keywords** Agent-based modeling and simulation · Levers of control · Replication · Imitation · *NKCS*-framework

## 1 Introduction

Multi-unit organizations are (potentially geographically) dispersed organizations that consist of a large number of divisions such as retail chain stores or fast-food franchises [9]. The divisions in multi-unit organizations predominantly operate in the same industry and promise customers the same brand experience in all divisions [18].

R. Khodzhimatov (✉)
Digital Age Research Center, University of Klagenfurt, Klagenfurt 9020, Austria
e-mail: ravshanbek.khodzhimatov@aau.at

S. Leitner · F. Wall
Department of Management Control and Strategic Management, University of Klagenfurt, Klagenfurt 9020, Austria
e-mail: stephan.leitner@aau.at

F. Wall
e-mail: friederike.wall@aau.at

To achieve this, the organizations have to make sure that divisions comply with the standards and best practices, but at the same time are given enough freedom to discover successful practices in the first place [2, 27].

To address this tension, organizations may employ different control mechanisms. Simons' Levers of Control framework [22] identifies four levers that constitute a management control system. Diagnostic control systems are formal mechanisms that ensure that the branches work towards the agreed-upon goal (e.g., incentive schemes). Interactive control systems are formal information systems which give a focused view on the aspects of performance (e.g., KPIs). Boundary systems delineate the acceptable behavior in the organization (e.g., codes of conduct, franchise operations manuals) and can be enabling or constraining, depending on the management's decisions. Belief systems is a set of core organizational values and definitions that management uses to foster a desired environment (e.g., mission statements, organizational culture).

Belief systems in this context can be used to describe *conformity*, which is defined as the internal desire of individuals to alter their behavior to match that of their peers. In contrast to compliance to organizational requirements, individuals conform voluntarily in pursuit of goals to blend into a team, gain approval of others, or increase accuracy of their actions by adopting the best practices of their peers [4]. The actual desire to conform changes with cultural and demographic characteristics of individuals and with the environment and norms [5]. However, the effect of desire of branch managers to conform on the actual adoption of the peers' practices depends on many factors, including the similarity in faced tasks [6], geographic proximity between branches [7], communication channels between branches [11], and the rotation of employees [13].

Chang and Harrington [3] studied the extent of centralization in retail chains in which managers come up with ideas for new practices. They found that organizations should employ boundary control systems that allow branches to experiment with new practices and routines subject to the constraint that branches need to adopt predetermined practices, and to combine this with diagnostic control systems by rewarding branch managers for replicating and passing along ideas. Garvin and Levesque [9] proposed diagnostic control systems that allow managers to prioritize their branch performance over the adoption of set practices. These studies implicitly assume that the best practices can be *codified* and put in a guideline, ready for replication. However, Haldin-Herrgard [10] showed that this is usually not possible, and suggested decentralized (interpersonal) knowledge transfer mechanisms. Additionally, Garvin and Levesque [9] found that due to the large number of branches, it is difficult to enforce a centralized control in multi-unit organizations. In this context, the less studied lever of control, belief systems, may be more effective because they allow organizations to foster an environment for imitation without strict centralized control mechanisms [23].

In this study we are interested to what extent does the individuals' desire to conform affect the diffusion of knowledge between divisions, and how do different network structures through which the agents communicate affect this relation. We employ an agent-based simulations approach [19, 26] to model the multi-unit

organization and the *NKCS*-framework [14, 15] to model the environments in which the units (divisions) operate. The rest of the paper is structured as follows: Sect. 2 presents the method, Sect. 3 summarizes our findings, Sect. 4 concludes the paper.

## 2 Model

In this section we introduce the agent-based model of an organization with $P = 5$ units. The task environment is based on the *NKCS*-framework [15, 20]. Agents make decisions to (a) increase their performance and (b) conform to the observed behavior of others. Section 2.1 introduces the task environment, Sects. 2.2 and 2.3 characterize the agents and describe how conformity is modeled. Section 2.4 describes the agents' search process, and Sect. 2.5 provides an overview of the sequence of events in the simulation.

### 2.1 Task Environment

We model an organization with $P = 5$ agents (units), each of which faces a complex decision problem that is expressed as the vector of $N = 4$ binary choices:

$$\mathbf{x} = (\underbrace{x_1, x_2, x_3, x_4}_{\mathbf{x}^1}, \ldots, \underbrace{x_{17}, x_{18}, x_{19}, x_{20}}_{\mathbf{x}^5}), \tag{1}$$

where bits $x_i \in \{0, 1\}$ represent single tasks. Every decision on a task $x_i$ yields a uniformly distributed performance contribution $\phi(x_i) \sim U(0, 1)$. The decision problem is *complex* in that the performance contribution $\phi(x_i)$, might be affected not only by the decision $x_i$, but also by decisions $x_j$, where $j \neq i$.

We differentiate between two types of such interdependencies: (a) *internal*, in which interdependence exists between the tasks within unit $p$, and (b) *external*, in which interdependence exists between the tasks in units $p$ and $q$ for $p \neq q$. We control interdependencies by parameters $K, C, S$, so that every task interacts with exactly $K$ other tasks internally and $C$ tasks assigned to $S$ other agents externally [14]:

$$\phi(x_i) = \phi(x_i, \underbrace{x_{i_1}, \ldots, x_{i_K}}_{\substack{K \text{ internal} \\ \text{interdependencies}}}, \underbrace{x_{i_{K+1}}, \ldots, x_{i_{K+C \cdot S}}}_{\substack{C \cdot S \text{ external} \\ \text{interdependencies}}}), \tag{2}$$

where $i_1, \ldots, i_{K+C \cdot S}$ are distinct and not equal to $i$. The exact choice of the coupled tasks is random with one condition: every task affects and is affected by exactly $K + C \cdot S$ other tasks. In our analysis we consider two benchmark cases: (i) only internal interdependence ($K = 3, C = S = 0$) and (ii) internal and external interdependence: ($K = C = S = 2$), as depicted in Fig. 1.

| | $\phi_1$ | $\phi_2$ | $\phi_3$ | $\phi_4$ | $\phi_5$ | $\phi_6$ | $\phi_7$ | $\phi_8$ | $\phi_9$ | $\phi_{10}$ | $\phi_{11}$ | $\phi_{12}$ | $\phi_{13}$ | $\phi_{14}$ | $\phi_{15}$ | $\phi_{16}$ | $\phi_{17}$ | $\phi_{18}$ | $\phi_{19}$ | $\phi_{20}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $x_1$ | x | x | x | x | | | | | | | | | | | | | | | | |
| $x_2$ | x | x | x | x | | | | | | | | | | | | | | | | |
| $x_3$ | x | x | x | x | | | | | | | | | | | | | | | | |
| $x_4$ | x | x | x | x | | | | | | | | | | | | | | | | |
| $x_5$ | | | | | x | x | x | x | | | | | | | | | | | | |
| $x_6$ | | | | | x | x | x | x | | | | | | | | | | | | |
| $x_7$ | | | | | x | x | x | x | | | | | | | | | | | | |
| $x_8$ | | | | | x | x | x | x | | | | | | | | | | | | |
| $x_9$ | | | | | | | | | x | x | x | x | | | | | | | | |
| $x_{10}$ | | | | | | | | | x | x | x | x | | | | | | | | |
| $x_{11}$ | | | | | | | | | x | x | x | x | | | | | | | | |
| $x_{12}$ | | | | | | | | | x | x | x | x | | | | | | | | |
| $x_{13}$ | | | | | | | | | | | | | x | x | x | x | | | | |
| $x_{14}$ | | | | | | | | | | | | | x | x | x | x | | | | |
| $x_{15}$ | | | | | | | | | | | | | x | x | x | x | | | | |
| $x_{16}$ | | | | | | | | | | | | | x | x | x | x | | | | |
| $x_{17}$ | | | | | | | | | | | | | | | | | x | x | x | x |
| $x_{18}$ | | | | | | | | | | | | | | | | | x | x | x | x |
| $x_{19}$ | | | | | | | | | | | | | | | | | x | x | x | x |
| $x_{20}$ | | | | | | | | | | | | | | | | | x | x | x | x |

(a) Internal interdependence     (b) External interdependence

**Fig. 1** Stylized interdependence structures with $P = 5$ agents facing $N = 4$ binary tasks. The crossed cells indicate inter-dependencies as follows: let $(i, j)$ be coordinates of a crossed cell in row-column order, then performance contribution $\phi(x_i)$ depends on decision $x_j$

Using Eq. 2, we generate *performance landscapes* as follows: for every task $x_i$ we generate performance contribution values corresponding to every combination of interdependent decisions from a uniform distribution. This results in a $N \times 2^{1+K+C \cdot S}$ matrix of uniform random numbers. We generate entire landscapes at the beginning of every simulation run to find the overall global maximum and normalize our results accordingly, to ensure comparability among different simulation runs.

At each time period $t$, agent $p$'s performance is a mean of performance contributions of tasks assigned to that agent:

$$\phi^P(\mathbf{x}_t^p) = \frac{1}{N} \sum_{x_i \in \mathbf{x}_t^p} \phi(x_i), \tag{3}$$

and the organization's overall performance is a mean of all agents' performances:

$$\Phi(\mathbf{x}_t) = \frac{1}{P} \sum_{p=1}^{P} \phi^P(\mathbf{x}_t^p). \tag{4}$$

All agents in the multi-unit organization operate in (not perfectly) similar environments, i.e., the same decisions of agents tend to lead to the same performance, with minor differences stemming from their local environments. We model this similarity between task environments of different units using the pairwise correlations between their performance landscapes [25]:

(a) Star network          (b) Ring network          (d) Cycle network



(e) Line network

**Fig. 2** Network structures in which $P = 5$ agents (nodes) share information through directed links. Star and Line networks capture unidirectional transfer of knowledge and Ring and Cycle networks capture mutual transfer of knowledge

$$\mathbf{corr}\left(\phi^p(\mathbf{x}_i^p), \phi^q(\mathbf{x}_i^q)\right) = \rho \in [0, 1], \tag{5}$$

for all $1 \leq i \leq N$ and $p \neq q$. In our analysis we use the value of $\rho = 0.9$ as our benchmark, as it represents high similarity in units with a few local differences.

## 2.2 Conformity Metric

To measure conformity we implement our version of the Social Cognitive Optimization algorithm introduced by Xie et al. [28]. At every time step $t$, agents share the decisions they have made on their tasks with the fellow agents, according to one of the network structures described in Fig. 2, where nodes represent agents and the directed links represent sharing of information.[1] Every agent $p$ stores the shared information in the memory set $L^p$ for up to $T_L = 50$ periods, after which the information is "forgotten" (removed from $L^p$).

The measure of conformity of agent $p$'s decisions $\mathbf{x}_t^p$ is computed as the average of the matching bits in the memory:

$$\phi_{\mathrm{conf}}(\mathbf{x}_t^p) = \begin{cases} \dfrac{1}{|L_t^p| \cdot N} \displaystyle\sum_{\mathbf{x}^L \in L_t^p} \sum_{i=1}^{N} [x_i^p == x_i^L], & t > T_L \\ 0, & t \leq T_L \end{cases} \tag{6}$$

---

[1] Including network structures in the model is a major extension over the papers using a similar approach to model the spread of information in organizations [16, 17].

where $|L_t^p|$ is the number of entries in agent $p$'s memory at time $t$, and the statement inside the square brackets is equal to 1 if true, and 0 if false [12].

## 2.3 Agents' Preferences

We model the agents' preferences as a linear function [1, 8, 24] of performance $\phi^p$ and conformity metric $\phi_{\text{conf}}$ (see Eqs. 3 and 6):

$$u^p(\mathbf{x}^p) = \alpha \cdot \phi^p(\mathbf{x}^p) + \beta \cdot \phi_{\text{conf}}(\mathbf{x}^p) \tag{7}$$

where $\alpha + \beta = 1$.

## 2.4 Search Process

In line with Simon [21], our agents are *boundedly rational*. In particular, the agents are not global optimizers and want to increase their utility given limited information: at time $t$, agents can observe their own performance in the last period, $\phi^p(\mathbf{x}_{t-1}^p)$, and the decisions of all team members in the last period *after* they are implemented, $\mathbf{x}_{t-1}$.

In order to come up with new solutions to their decision problems, agents perform a search in the neighbourhood of $\mathbf{x}_{t-1}$ as follows: agent $p$ randomly switches one decision $x_i \in \mathbf{x}^p$ (from 0 to 1, or vice versa), and assumes that other agents will not switch their decisions (Levinthal [20] describes situations in which agents switch more than one decision at a time as *long jumps* and states that such scenarios are less likely to occur, as it is hard or risky to change multiple processes simultaneously). We denote this vector with one switched element by $\hat{\mathbf{x}}_t^p$.

Next, the agent has to make a decision whether to stick with the status quo, $\mathbf{x}_t^p$, or to switch to the newly discovered $\hat{\mathbf{x}}_t^p$. The rule for this decision is to maximize the utility function defined in Eq. 7:

$$\mathbf{x}_t^p = \underset{\mathbf{x} \in \{\mathbf{x}_{t-1}^p, \hat{\mathbf{x}}_t^p\}}{\arg\max}\, u(\mathbf{x}), \tag{8}$$

## 2.5 Process Overview, Scheduling and Main Parameters

The simulation model has been implemented in *Python 3.8* and *Numba* just-in-time compiler. Every simulation round starts with the initialization of the agents' performance landscapes, the assignment of tasks to $P = 5$ agents. For reliable results, we generate the entire landscapes before the simulation, which is feasible for $P = 5$ given modern computing limitations, and the initialization of an $M = 20$ dimensional

**Fig. 3** Process overview

**Table 1** Main parameters

| Parameter | Description | Value |
|---|---|---|
| $M$ | Total number of tasks | 20 |
| $P$ | Number of agents | 5 |
| $N$ | Number of tasks assigned to a single agent | 4 |
| [K,C,S] | Internal and external couplings | [3, 0, 0], [2, 2, 2] |
| $\rho$ | Pairwise correlation between landscapes | 0.9 |
| $T_L$ | Memory span of agents | 50 |
| $T$ | Observation period | 500 |
| $R$ | Number of simulation runs per scenario | 1000 |
| $[\alpha, \beta]$ | Weights for performance $\phi$ and conformity $\phi_{\text{conf}}$ | [1, 0], [0.5, 0.5], [0, 1] |

bitstring as a starting point of the simulation run. After initialization, agents perform the *hill climbing* search procedure outlined above and share information regarding their own decisions according to the network structure. The observation period $T$, the memory span of the employees $T_L$, and the number of repetitions in a simulation, $R$, are exogenous parameters, whereby the latter is fixed on the basis of the coefficient of variation. Figure 3 provides an overview of this process and Table 1 summarizes the main parameters used in this paper.

# 3   Results

In this section we present selected findings from running $R = 1000$ simulations for 4 different network structures and task environments with and without external interdependencies between departments. In each simulation scenario we observe organization-level performance and the measure of synchrony across divisions for $T = 500$ time periods. Section 3.1 defines the synchrony measure, and Sects. 3.2 and 3.3 present the findings.

## 3.1   Measure of Synchrony

To measure the synchrony of a strategy we first define the Hamming distance, which is a metric that returns the number of bits that are distinct in two bit strings. For example, the Hamming distance between a bit string 1001 and 1101 is equal to 1, as they differ in only one bit, and flipping just one bit is sufficient to make them equal. Similarly, the Hamming distance between identical bit strings 1001 and 1001 is equal to zero.

Next, we define the *asynchrony* or distinctness of a bitstring as the sum of all pairwise Hamming distances between the bit sub-strings allocated to different agents.

$$\mathbf{H}(\mathbf{x}) = \sum_{p=1}^{P} \sum_{q=p}^{P} H(\mathbf{x}^p, \mathbf{x}^q) \tag{9}$$

Finally, we measure the synchrony of a bitstring as a complement of the asynchrony normalized by its maximum:

$$\mathbf{S}(\mathbf{x}) = 1 - \frac{\mathbf{H}(\mathbf{x})}{\max\{\mathbf{H}(\mathbf{x})\}} \tag{10}$$

## 3.2   Findings Regarding Synchrony

In this section we analyze how the synchrony in the organization is affected by the individuals' desire to conform by the different network structures through which they communicate. Solely prefer the conformity and do not consider the actual performance in their departments (i.e., $\alpha = 0$, $\beta = 1$) to understand how the network structures operate to spread the best practices between departments. Figure 4 shows that the top-to-bottom unilateral sharing of knowledge in Star and Line network structures leads to the full synchrony, with Line network being slower to converge. The peer-to-peer egalitarian network structures like Cycle network lead only to par-

**Fig. 4** Synchrony measure for full conformity ($\alpha = 0$, $\beta = 1$). This figure applies to all scenarios with $P = 5$, and, by construction, is not affected by the structure of the task environment

tial synchrony even if the agents have a full desire to conform. By construction, this relation holds for all task environments with $P = 5$ agents.

This intuitive finding, however, no longer holds in a more realistic scenario, in which agents have both performance and conformity in their preferences (i.e. $\alpha = \beta = 0.5$). Indeed, we find that centralized Star network leads to a high synchrony only in the short term in the absence of external interdependencies between divisions. In the long term, however, the Cycle network leads to a higher synchrony. In presence of external interdependencies between divisions, the Ring network leads to the highest synchrony and passes the Star network after 50 periods. Moreover, the Line network leads to the lowest synchrony for interrelated divisions. All of these scenarios, however, lead to a significantly higher synchrony than the situations in which agents do not have a desire to conform (i.e., $\alpha = 1$, $\beta = 0$), which lead to a synchrony measure $\mathbf{S}(\mathbf{x}) \leq 0.3$ (Fig. 5).

These results indicate that, while conformity can significantly increase the diffusion of best practices in the organizations, the management should be careful in promoting it and consider the nature of tasks the units are facing and the interdependence between them. We find that the naive idea of identifying the successful unit and promoting other units to directly replicate its practices does not always lead to the highest diffusion of knowledge, and that forsaking a centralized control and promoting a peer-to-peer communication is more effective in the long run.

### 3.3 Performance Measure

Next, we look at how the different network structures affect the organization-level performance for environments with and without external interdependencies. We find

(a) Internal interdependence    (b) External interdependence

**Fig. 5** Synchrony measure for moderate level of conformity



(a) Internal interdependence    (b) External interdependence

**Fig. 6** Performance for moderate level of conformity

that the centralized Line and Star network structures actually lead to less organizational performance than the decentralized Ring and Cycle networks. This happens because in the centralized network structures, the central units do not see decisions of their peers and, thus, cannot benefit from the knowledge they have gained. In the decentralized network structures, on the other hand, all units directly or indirectly observe the decisions made by their peers, and, thus, can benefit from the knowledge of all departments.

Between the latter two, the Ring network leads to the higher performance in presence of external interdependencies between units, and the Cycle network leads to the highest long-term performance in the absence of them (Fig. 6).

## 4 Conclusion

In this paper we studied how multi-unit organizations can employ individuals' belief systems, particularly, their desire to conform to the behavior of their peers, to achieve the diffusion of best practices between their units. We performed an agent-based simulation of the organization and compared the achieved synchrony for different network structures between the units. We found that, contrary to the intuition, the centralized spread of knowledge from a single unit to others leads to a lower long-term synchrony than the decentralized peer-to-peer sharing of knowledge between all units. Interestingly, our results do not feature a trade-off between organizational performance and synchrony—we find that in most situations the decentralized Cycle and Ring networks help to achieve both the high synchrony and the high performance.

The implications of our results are that the management in multi-unit organizations should forsake the centralized control over the diffusion of knowledge (via diagnostic or boundary control systems) and to promote an organizational culture of sharing knowledge and conforming to the most frequent practices. The limitations of our research include the lack of historical performance in the agents' consideration to conform—further research might address this via dynamically updated weights for the desire to conform.

## References

1. Akerlof, G.A.: A theory of social custom, of which unemployment may be one consequence. Q. J. Econ. **94**(4), 749 (1980)
2. Argote, L., Ingram, P.: Knowledge transfer: a basis for competitive advantage in firms. Organ. Behav. Hum. Decis. Process. **82**(1), 150–169 (2000)
3. Chang, M.H., Harrington, J.E.: Centralization vs. decentralization in a multi-unit organization: a computational model of a retail chain as a multi-agent adaptive system. Manage. Sci. **46**(11), 1427–1440 (2000)
4. Cialdini, R.B., Goldstein, N.J.: Social influence: compliance and conformity. Annu. Rev. Psychol. **55**(1), 591–621 (2004)
5. Cialdini, R.B., Reno, R., Kallgren, C.: A focus theory of normative conduct: recycling the concept of norms to reduce littering in public places. J. Pers. Soc. Psychol. **58**, 1015–1026 (06 1990)
6. Darr, E.D., Argote, L., Epple, D.: The acquisition, transfer, and depreciation of knowledge in service organizations: productivity in franchises. Manage. Sci. **41**(11), 1750–1762 (1995)
7. Galbraith, C.S.: Transferring core manufacturing technologies in high-technology firms. Calif. Manage. Rev. **32**(4), 56–70 (1990)
8. Gali, J.: Keeping up with the joneses: consumption externalities, portfolio choice, and asset prices. J. Money Credit Bank. **26**(1), 1–8 (1994)
9. Garvin, D.A., Levesque, L.C.: The multiunit enterprise. Harvard Bus. Rev. 106–107 (2008)
10. Haldin-Herrgard, T.: Difficulties in diffusion of tacit knowledge in organizations. J. Intellect. Capital **1**(4), 357–365 (2000)
11. Hansen, M.T.: Knowledge networks: explaining effective knowledge sharing in multiunit companies. Organ. Sci. **13**(3), 232–248 (2002)
12. Iversion, K.E.: A programming language. Wiley (1962)

13. Kane, A.A., Argote, L., Levine, J.M.: Knowledge transfer between groups via personnel rotation: effects of social identity and knowledge quality. Organ. Behav. Hum. Decis. Process. **96**(1), 56–71 (2005)
14. Kauffman, S.A., Johnsen, S.: Coevolution to the edge of chaos: coupled fitness landscapes, poised states, and coevolutionary avalanches. J. Theor. Biol. **149**(4), 467–505 (1991)
15. Kauffman, S.A., Weinberger, E.D.: The NK model of rugged fitness landscapes and its application to maturation of the immune response. J. Theor. Biol. **141**, 211–245 (1989)
16. Khodzhimatov, R., Leitner, S., Wall, F.: Interactions between social norms and incentive mechanisms in organizations. In: Coordination, Organizations, Institutions, Norms, and Ethics for Governance of Multi-agent Systems (2021)
17. Khodzhimatov, R., Leitner, S., Wall, F.: On the effect of social norms on performance in teams with distributed decision makers. In: Social, Cultural, and Behavioral Modeling, pp. 299–309. Springer International Publishing, Cham (2021)
18. Kim, S.: The rise of multiunit firms in U.S. manufacturing. Explor. Econ. Hist. **36**(4), 360–386 (1999)
19. Leitner, S., Wall, F.: Simulation-based research in management accounting and control: an illustrative overview. J. Manage. Control **26**(2–3), 105–129 (2015)
20. Levinthal, D.A.: Adaptation on rugged landscapes. Manage. Sci. **43**(7), 934–950 (1997)
21. Simon, H.A.: Models of Man; Social and Rational (1957)
22. Simons, R.: Control in an age of empowerment. Harvard Bus. Rev. (1995)
23. Tessier, S., Otley, D.: A conceptual development of Simons' levers of control framework. Manage. Acc. Res. **23**(3), 171–185 (2012)
24. Tversky, A., Kahneman, D.: Loss aversion in riskless choice: a reference-dependent model. Q. J. Econ. **106**(4), 1039–1061 (1991). https://doi.org/10.2307/2937956
25. Verel, S., Liefooghe, A., Jourdan, L., Dhaenens, C.: On the structure of multiobjective combinatorial search space: MNK-landscapes with correlated objectives. Eur. J. Oper. Res. **227**(2), 331–342 (2013)
26. Wall, F., Leitner, S.: Agent-based computational economics in management accounting research: opportunities and difficulties. J. Manage. Acc. Res. **33**(3), 189–212 (11 2020)
27. Winter, S.G., Szulanski, G.: Replication as strategy. Organ. Sci. **12**(6), 730–743 (2001)
28. Xie, X.F., Zhang, W.J., Yang, Z.L.: Social cognitive optimization for nonlinear programming problems. In: Proceedings of International Conference on Machine Learning and Cybernetics, vol. 2, pp. 779–783. IEEE (2002)

# Combining Experiments with Agent-Based Modeling: Benefits for Experimental Management Accounting Research

**Jannick Plähn** ⓘ**, Lucia Bellora-Bienengräber** ⓘ**, Kai G. Mertens** ⓘ**, and Matthias Meyer** ⓘ

**Abstract** Laboratory experiments are among the most frequently used methods in management accounting research because they offer high internal validity, enabling the examination of causal relationships. However, experiments often struggle with providing support for a specific proposed causal mechanism, given the abundance of psychological and behavioral theories that predict similar outcomes. In this paper, we argue that agent-based modeling is well suited to complement experiments because agent-based modeling is a powerful method to increase confidence in the proposed causal mechanism. As a showcase project, we conduct an experiment to explain antecedents of honest reporting behavior in a participative budgeting setting and propose that a social norm of honesty is the underlying causal mechanism. Next, we adapt an agent-based model to our participative budgeting setting and create two submodels incorporating alternative causal mechanisms. Finally, we assess the capability of the two submodels to reproduce the experiment's results to evaluate whether the observed behavior in the experiment can be better explained with the causal mechanism representing social norm theory.

**Keywords** Participative budgeting · Social norm · Mixed-method · Experiment · Agent-based modeling

J. Plähn (✉) · K. G. Mertens · M. Meyer
Institute of Management Accounting and Simulation, Hamburg University of Technology, Am Schwarzenberg-Campus 4, 21073 Hamburg, Germany
e-mail: jannick.plaehn@tuhh.de

L. Bellora-Bienengräber
Department of Accounting and Auditing, Faculty of Economics and Business, University of Groningen, Nettelbosje 2, 9747, AE Groningen, The Netherlands

Wilbur O. and Ann Powers College of Business, School of Accountancy, Clemson S.C. 29634, USA

# 1   Introduction

As in many other social science disciplines, laboratory experiments are among the most frequently used methods in management accounting research. Bloomfield, et al. [1] examined all papers published in the Journal of Accounting and Economics, the Journal of Accounting Research, the Journal of Accounting and Economics, The Accounting Review, and Accounting, Organizations and Society from 2003 to 2013 and found that experiments and field studies are the most used method in management accounting research to gather data. Guffey and Harp [2] examine articles published in the Journal of Management Accounting Research from 1989 to 2013 and find that experiments are among the five most frequently used methods. Besides the need for management accounting research to gather data themselves because available archival data is limited in this research field [1], experiments are often used because randomization and a controlled setting provide high internal validity, which allows for studying causal relationships [3].

However, experiments often struggle when researchers are interested in the intervening mechanisms through which the independent variable affects the dependent variable. Since experiments only provide direct evidence about individuals' behavior, but not their reasoning [4], it may be challenging to provide support for a specific mechanism due to the abundance of psychological mechanisms predicting similar outcomes [5–7]. Distinguishing between possible mechanisms may be crucial when the research goal is to derive effective interventions to change a certain behavior. For example, individuals may split their endowment in a dictator game equally because they have a social preference for fairness and value an equal distribution (i.e., they choose a fair behavior because they generally prefer fairness in social interactions) [8], or because they comply with an established social norm of fairness (i.e., they choose a fair behavior because they comply with their expectations about others' behavior and beliefs) [9]. Interventions that aim to increase fairness should be aimed at changing individuals' expectations about others' behavior and beliefs when a social norm is the underlying causal mechanism, while this is ineffective in the case of a mechanism rooted in the social preference for fairness.

In this regard, process evidence becomes crucial. Process evidence is data describing the underlying causal mechanisms through which independent variables affect dependent variables [6]. In addition to mediation and moderation analyses, Asay, et al. [6] propose a multiple-method approach to take advantage of triangulation and increase confidence in the proposed causal mechanism when there are multiple plausible mechanisms. For their multiple-method approach, Asay, et al. [6] consider experiments in combination with analytical models, archival analyses, surveys, and interviews. We extend the list of methods and argue that the combination of experiments with agent-based modeling (ABM) is also able to increase confidence in the proposed causal mechanism [4]. Although there are examples of combining experiments and ABM in other disciplines like behavioral economics (e.g., [10, 11]) or team cognition (e.g., [12, 13]), ABM as a research method is fairly new to management accounting [14]. Considering the frequent use of experiments

in management accounting research and the limitation of experiments to support a specific causal mechanism, we see great potential for ABM to benefit management accounting research. The following steps outline our approach to combine experiments and ABM:

1. Experimental research provides empirical information on how individuals behave in a particular situation.
2. An agent-based model is developed in which the proposed and alternative causal mechanisms are explicitly implemented as part of agents' decision process.
3. The results of the agent-based models based on the proposed and alternative causal mechanisms and the experiment are compared to check whether they resemble the experiment's results; if this is the case, this would provide support that the implemented mechanism proxies the participants' decision process in the experiment.

## 2 Experiment

### 2.1 Design

As an illustration of the advantages of combining experiments and ABM in experimental management accounting research, we examine honest reporting in the context of participative budgeting. We choose this setting for two reasons. First, participative budgeting is one of the most widely investigated topics of experimental research in management accounting [15]. Second, management accounting literature has recently begun to use social norm theory to predict and decrease opportunistic behavior (i.e., dishonest reporting) in participative budgeting [16–18]. However, compliance with a social norm can only be inferred if the behavior is shown to be conditional on individuals' expectations about others' behavior and beliefs; thus, observing behavior in an experiment (e.g., more or less honest reporting behavior) is not sufficient to infer the existence and influence of a social norm [19]. Therefore, we argue that the case of honest reporting in participative budgeting is an appropriate setting where ABM can increase confidence in the proposed mechanism (in our case, a social norm of honesty).

Our experiment employs a 2 (internal reporting environment closed/open) × 2 (pooled profit-sharing plan absent/present) × 10 (periods) mixed factorial design. Each participant acts as a division manager in an organization consisting of three divisions. In each period, each division yields revenues of 6000 Lira and pays participants a fixed salary of 500 Lira. Participants receive information about their division's true costs and must report their true costs to corporate headquarters to get funding. In line with the trust contract from Evans, et al. [24], corporate headquarters only knows that divisions have a possible range of costs between 4000 and 5500 Lira and accepts any budget request within this range. Thus, participants can increase their payoff by overstating true costs, while this decreases the division's and, subsequently, the

organization's profit. For example, when a participant chooses to report the highest possible costs of 5500 Lira, the division's profit for that period will be 0 Lira due to the participant's fixed salary of 500 Lira, and the division's contribution to the organization's profit will thus also be 0 Lira.

We manipulate the openness of the internal reporting environment by varying the information participants receive about others' true costs and their submitted budget reports to corporate headquarters. In the closed reporting environment, participants do not get information about others' true costs and submitted budget reports. In the open reporting environment, participants observe others' true costs and submitted budget reports; thus, participants can infer others' level of honesty. Further, we manipulate the presence of a pooled profit-sharing plan. When a pooled profit-sharing plan is absent, participants only receive a fixed salary of 500 Lira. When a pooled profit-sharing plan is present, participants receive their fixed salary plus 5% of the pooled profits of all three divisions in their organization (i.e., the organization's profit). In line with Boster et al. [20], we design the pooled profit-sharing plan in such a way that dishonest reporting still maximizes the participants' payoffs. Finally, participants in all conditions can sanction other participants in each period. In line with experiments in behavioral economics [21–23], we operationalize peer sanctioning as participants' option to assign an integer amount of 0–10 sanction points to each of the other two participants in their organization in each period. Each sanctioning point assigned decreases the payoff of the sanctioned participant by 10% with a maximum of 100% but is also costly for the sanctioning participant as each sanction point assigned decreases the sanctioning participant's payoff. For example, when participants assign eight sanction points in total to the other two participants, their costs for assigning sanction points equal their fixed salary. Our dependent variable is participant's honesty. In line with Evans, et al. [24], we measure honesty ranging from 0 to 1 as follows:

$$\text{Honesty} = 1 - (\text{Reported costs} - \text{True costs})/(5500 - \text{True costs}) \qquad (1)$$

Our predictions regarding participants' honesty are based on social norm theory [9]. It states that one reason for complying with a social norm (in our case, the social norm of honesty) is the expectation that relevant others (in our case, other division managers) follow and think one should follow a given behavioral rule (in our case, report honestly) and are willing to sanction otherwise. First, we predict that an open internal reporting environment will increase honesty by increasing division managers' perceived risk of being sanctioned, compared to a closed internal reporting environment. While division managers in the closed internal reporting environment can hide their misreporting behind the lack of transparency among division managers, their misreporting in an open internal reporting environment is exposed to their peers and may trigger peers' sanctions. Second, we predict that a pooled profit-sharing plan will further increase honesty. Since overstating costs in the budgeting process impacts the division's profit, the organization's profits, and subsequently, division managers' profit-share, pooled profit-sharing plans introduce interdependency among managers. This interdependency will increase the number

of sanctions because sanctions are driven by negative emotions like the feeling of being exploited by others [23, 25]. Being sanctioned by others will also increase the salience of the norm, increasing the psychological costs for future norm violations [26, 27].

## 2.2 Results

The average level of honesty per period and condition is shown in Fig. 1. To test our predictions, we use a mixed-effects regression with random effects at the individual level to account for within-subject dependency. We use participants' honesty in each period as the dependent variable and the openness of the internal reporting environment and the presence of a pooled profit-sharing plan as independent variables. We also control for potential time effects by including period as a covariate. Untabulated results of the mixed-effects regressions show that honesty increases by approximately 0.18 when the internal reporting environment changes from closed to open ($\beta = 0.18$, t-value = 2.504, $p < 0.05$). Similarly, honesty increases by approximately 0.14 when a pooled profit-sharing plan is present compared to when a pooled profit-sharing plan is absent ($\beta = 0.14$, t-value = 2.082, $p < 0.05$). Both effects are significant and in line with our prediction that an open internal reporting environment and a pooled profit-sharing plan increase honesty. In contrast, our results show a significant negative effect of period ($\beta = -0.02$, t-value = $-8.097$, $p < 0.01$), suggesting that honesty declines over time.

Although dishonest reporting always maximizes participants' payoffs in all conditions, our results show the highest level of honesty in Condition 4 where participants can see each other's reporting behavior and a pooled profit-sharing plan introduces interdependency among participants in terms of their payoff. Compared to the other conditions, we argue that in this condition the social norm of honesty is more salient (i.e., participants have higher empirical and normative expectations), which motivates more participants to comply with the social norm and report honestly. The reason is that in this condition participants' dishonest reporting decreases others' payoff due to the pooled-profit sharing plan. Since dishonest reporting can be observed due to the open internal reporting environment, we argue that among all conditions this condition has the highest number of sanctions. In line with our argumentation, untabulated results show that in Condition 4 dishonest participants are more often sanctioned than in the other conditions. Being sanctioned gives participants a clear signal that peers view dishonest reporting as a norm violation, thus increasing participants' empirical and normative expectations regarding the social norm of honesty [10]. However, an alternative explanation may be rooted in agency theory. Standard agency theory predicts that individuals are solely motivated by material self-interest and want to maximize their utility through wealth [28]. Since sanctions significantly decrease participants' payoff, reporting honestly and thereby avoiding any sanctions may become participants' reporting choice with the highest payoff. Thus, participants are still solely motivated by material self-interest and simply report more honestly

**Fig. 1** Average honesty per period and condition

to avoid sanctions that are more common in Condition 4. This alternative explanation only considers the material consequences of sanctions but not sanctions' norm-signaling function. To distinguish between these two explanations, we apply an agent-based model. In Sect. 3, we develop an agent-based model with submodels that represent both causal mechanisms. In Sect. 4, we assess whether the causal mechanisms proposed in social norm theory or agency theory can better approximate the observed honesty level in the experiment. Finally, based on this comparison, we draw conclusions about which mechanism is most likely underlying participants' behavior in the experiment.

## 3 Agent-Based Model

### 3.1 Setting, Procedure, and Agents' Decision Making

We adapt the agent-based model by Andrighetto, et al. [10] to our participative budgeting setting in Condition 4 (i.e., open internal reporting environment and pooled profit-sharing plan) as it explicitly incorporates compliance with a social norm as part of agents' decision process. Their model is dynamic in that agents observe other agents' behavior, changing their preference to follow a social norm over time and thus going beyond a purely static social preference model [10]. In our model, like in our experiment, each agent represents a division manager, and three managers form an organization. In each period, agents learn the true costs of their division and decide whether or not to report honestly and to sanction the other two agents for their reporting choice. To simplify the model, agents' reporting choice is binary: report the true costs (i.e., report honestly) or report the highest possible costs of 5500 Lira (i.e., report dishonestly). Agents who reported honestly in a given period can sanction agents who reported dishonestly after observing their reporting choices. In line with the experiment, agents can assign up to 10 sanction points that reduce the sanctioned agent's payoff by 10% per sanction point but are also costly for the sanctioning agent. The choice to report honestly depends on a probability that is updated every period as a function of agents' individual drive (ID) and normative drive (ND).

The ID reflects agents' goal to maximize their payoff not considering what the norm prescribes and is updated with a winner-stay-losers-change algorithm [29]. In the model, agents calculate their potential payoff separately for each reporting choice considering the true costs in this period and assuming the other two agents in their organization would report like in the previous period. If agents have received sanction points in a previous period, they consider these when calculating the potential payoff for reporting dishonestly. Then the ID moves towards the action that returns the potential higher payoff. The ND models agents' motivation to comply with the social norm, dependent on the norm's salience. Norm salience is an agent's perception regarding the importance of the social norm within the group and is updated every period for each agent according to the norm cues agents receive. Norm cues include own norm-compliance and norm-violation, observed norm-compliance and norm-violation of other agents in the same organization, and being sanctioned by others. The effect of these norm cues on the norm salience varies and is derived from Cialdini, et al. [27] and defined in Andrighetto, et al. [10] (for details, see [10]). Since agents may receive different norm cues to update norm salience, there is heterogeneity regarding norm salience within the population. The extent to which ID and ND affect the reporting choice is determined through the parameters individual weight (IW) and normative weight (NW). These parameters express the importance of ID and ND for each agent. Using agent's ID and ND, as well as the parameters IW and NW, agent's probability to report honestly $p$ is calculated as follows:

$$\text{p}^{\text{honest in t}} = \text{p}^{\text{honest in t}-1} + (\text{ID} \times \text{IW} + \text{ND} \times \text{NW}) \tag{2}$$

In sum, in each period each agent goes through the following steps:

1. Process information from the previous period (starts in Period 2)

   (a) Update ID
   (b) Update ND
   (c) Update the probability to report honestly

2. Choose to report honestly (true costs of the period) or dishonestly (highest possible costs of 5500 Lira)
3. Observe the other two agents' reporting behavior
4. Choose to assign sanction points to the other two agents.

## 3.2 Submodels

The parameters IW and NW express the importance of ID and ND for each agent and always add up to 1 (i.e., NW = 1−IW), which is supposed to represent that a higher importance of material payoffs decreases the importance of complying with the social norm and vice versa. We vary the weights of IW and NW to build two different submodels with theory-compliant agents:

- Submodel 1: Agents based on agency theory (IW = 1.0, NW = 0.0)
- Submodel 2: Agents based on social norm theory (IW = 0.5, NW = 0.5)

In Submodel 1, agents' IW is 1 and NW is 0, thus agents only consider the material payoff of an action when updating their probability to report honestly in the next period. Since agency theory assumes that the maximization of wealth is agents' sole motivation [28], we consider this a good approximation for an agent based on agency theory. In Submodel 2, agents' IW is 0.5 and NW is 0.5, thus agents consider their wealth and compliance with the social norm of honesty when updating their probability to report honestly in the next period. We do not use agents who solely consider norm compliance (i.e., IW = 0, NW = 1) when updating their probability to report honestly in the next period since social norm theory postulates a norm-based utility function in which individuals tradeoff utility from material possessions against disutility from norm violation [9]. Therefore, we consider agents with an IW of 0.5 and an NW of 0.5 as a good approximation for an agent based on social norm theory.

For each submodel, we conduct 100 simulation runs in which 30 agents each are created and divided into groups of three to form 10 organizations as in the experiment. Further, each simulation run was performed with 10 periods representing the 10 budgeting periods. In each period in each run, we calculate the honesty level of the population, which is calculated as the average honesty of all 30 agents. Therefore, we use the same true costs as in the experiment. Further, we determine agents' initial probability to report honestly and to assign sanction points using the empirical data from the first period of the experiment. In the subsequent periods, the probability of sanctioning a dishonest agent is inversely proportional to the number of dishonest agents within the organization [10].

## 3.3 Results

Figure 2 shows the average honesty of the 100 populations in each period for both submodels. We can see that in Submodel 1 where agents solely consider their material payoffs, average honesty rapidly declines after Period 1 and remains at a very low level throughout the remaining nine periods. Although the level of honesty is very low, it is above 0 because in some organizations agents are sanctioned and therefore honest reporting becomes the wealth-maximizing reporting choice. Average honesty in Submodel 2 where agents balance wealth maximization and norm compliance also declines over time but this process is much slower than in Submodel 1. After Period 6, average honesty remains constant on a low level but is slightly higher than in Submodel 1. The results indicate that the social norm influences agents' behavior decreases over time due to norm violations, thus, agents consider norm compliance less and less over time, making considerations of wealth maximization relatively more important.
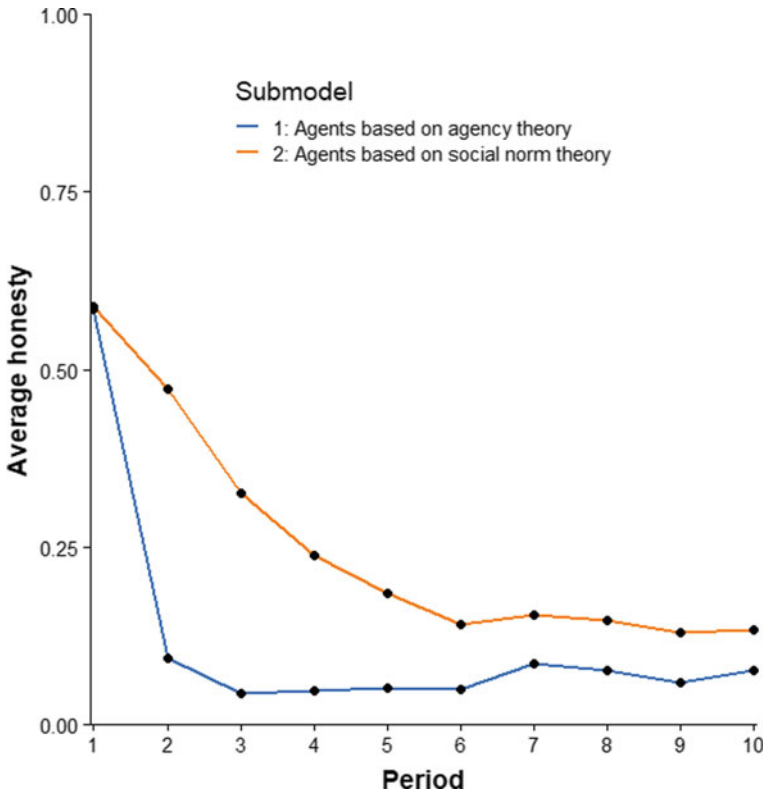


**Fig. 2** Average honesty per period and submodel

## 4    Comparison of Simulated and Empirical Honesty Levels

We compare simulated and empirical honesty levels to evaluate which causal mechanism is a better explanation for the experimental results. Since we predicted the strongest effect of a social norm of honesty in the condition with an open internal reporting environment and a pooled profit-sharing plan, we adapted the model to this situation and therefore use the average honesty of the participants from condition 4 as an empirical benchmark for our simulated honesty levels. We follow Lorscheid and Meyer [12] and calculate for both submodels the empirical distance as the absolute distance from the simulated honesty level to the empirical honesty level in each period $t$.

$$\text{Empirical distance}_t = \left| \text{simulated honesty}_t - \text{empirical honesty}_t \right| \qquad (3)$$

Table 1 shows the empirical distance averaged over all 100 runs per period and in total for all periods. For example, 0.08 in Period 1 from Submodel 1 expresses that in Period 1 the absolute honesty difference between the simulated population and participants from Condition 4 in the experiment is on average 0.08 (honesty ranges from 0 to 1).

The empirical distance averaged over all 10 periods is significantly lower in Submodel 2 (all = 0.22) than in Submodel 1 (all = 0.51). Thus, agents based on social norm theory provide honesty levels more consistent with the empirical honesty levels than agents based on agency theory. Further, the empirical distance of Submodel 1 is very high. Since considering the monetary consequences of sanctioning is the only reason for an agent in Submodel 1 to report honestly, this indicates that participants in the experiment not only report honestly to avoid the monetary consequences of being sanctioned. Submodel 2 shows significantly lower empirical distances than Submodel 1, which suggests that participants' behavior in the experiment can be better explained when considering the norm-signaling function of sanctions and participants' motivation to comply with a salient social norm of honesty. Nevertheless, also Submodel 2 still shows significant empirical distances, thus, there are other influences that Submodel 2 does not capture.

**Table 1**  Average empirical distances per period and submodel compared to Condition 4 in the experiment

| Submodel | Period | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | All |
| 1 | 0.08 | 0.67 | 0.61 | 0.36 | 0.47 | 0.58 | 0.48 | 0.66 | 0.63 | 0.58 | 0.51 |
| 2 | 0.07 | 0.18 | 0.14 | 0.29 | 0.29 | 0.23 | 0.36 | 0.19 | 0.22 | 0.28 | 0.22 |

# 5    Conclusion

Experiments are among the most frequently used methods in management accounting research because they allow us to examine causal relationships due to their high internal validity. However, experiments often struggle with providing support for a specific proposed causal mechanism, given the abundance of psychological and behavioral theories that predict similar outcomes. Adding to the proposed multi-method approach by Asay, et al. [6] to increase confidence in a specific causal mechanism, we suggest that ABM can supplement experiments by providing additional support for the proposed mechanism (in our case, a social norm of honesty).

To illustrate this, we adapted an agent-based model to our participative budgeting setting with an open internal reporting environment and a pooled profit-sharing plan. We created two submodels in which agents' decision process incorporates either the causal mechanism based on social norm theory or agency theory. Results show that agents who consider the monetary consequences and norm-signaling function of sanctions provide honesty levels more consistent with the empirical honesty levels than agents who in line with agency theory only consider the monetary consequences of sanctions. This provides initial support that participants in the experiment not only report honestly because they want to avoid the costs of being sanctioned but also because peers' sanctions increase participants' empirical and normative expectations that motivate more participants to comply with the social norm and report honestly. Since honesty levels in Submodel 2 are still different from the honesty levels in the experiment, this is only the first step to providing support for a social norm as the underlying causal mechanism in our experiment. Future research should refine our model. So far, Submodel 2 only considers heterogeneity regarding the norm information agents receive, but not regarding the extent to which agents are affected by the social norm. Since social norm theory assumes that individuals differ regarding the extent to which social norms affect their behavior [9], this appears to be a promising way to refine the model and test whether this further decreases empirical distances. Overall, this study provides a first step to support our claim that ABM is suitable to be part of a multi-method approach to increase confidence in a specific causal mechanism. We hope to pave the ground for further research questions in management accounting research to be tackled with such a multi-method approach.

# References

1. Bloomfield, R., Nelson, M.W., Soltes, E.: Gathering data for archival, field, survey, and experimental accounting research. J. Account. Res. **54**(2), 341–395 (2016)
2. Guffey, D.M., Harp, N.L.: The journal of management accounting research: a content and citation analysis of the first 25 years. J. Manag. Account. Res. **29**(3), 93–110 (2016)
3. Abernethy, M.A., Chua, W., Luckett, P.F., Selto, F.H.: Research in managerial accounting: learning from others' experiences. Account. Finan. **39**(1), 1–27 (1999)
4. Janssen, M.A., Ostrom, E.: Empirically based, agent-based models. Ecol. Soc. **11**(2) (2006)
5. Levitt, S.D., List, J.A.: Homo economicus evolves. Science **319**(5865), 909–910 (2008)

6.  Asay, H.S., Guggenmos, R., Kadous, K., Koonce, L., Libby, R.: Theory testing and process evidence in accounting experiments. Account. Rev., (2021)
7.  Smith, E.B., Rand, W.: Simulating macro-level effects from micro-level observations. Manage. Sci. **64**(11), 5405–5421 (2018)
8.  Fehr, E., Schmidt, K.M.: A theory of fairness, competition, and cooperation. Q. J. Econ. **114**(3), 817–868 (1999)
9.  Bicchieri, C.: The Grammar of Society: The Nature and Dynamics of Social Norms. Cambridge University Press, New York (2006)
10. Andrighetto, G., Brandts, J., Conte, R., Sabater-Mir, J., Solaz, H., Villatoro, D.: Punish and voice: punishment enhances cooperation when combined with norm-signalling. Plos One, **8**(6) (2013)
11. Klingert, F.M.A., Meyer, M.: Effectively combining experimental economics and multi-agent simulation: suggestions for a procedural integration with an example from prediction markets research. Comput. Math. Organ. Theory **18**(1), 63–90 (2012)
12. Lorscheid, I., Meyer, M.: Toward a better understanding of team decision processes: combining laboratory experiments with agent-based modeling. J. Bus. Econ. **91**, 1–37 (2021)
13. Lorscheid, I., Meyer, M.: Agent-based mechanism design—investigating bounded rationality concepts in a budgeting context. Team Perform. Manage. Int. J. **23**(1/2), 13–27 (2017)
14. Wall, F., Leitner, S.: Agent-based computational economics in management accounting research: opportunities and difficulties. J. Manag. Account. Res. **33**(3), 189–212 (2021)
15. Brown, J.L., Evans, J.H., III., Moser, D.V.: Agency theory and participative budgeting experiments. J. Manag. Account. Res. **21**(1), 317–345 (2009)
16. Douthit, J.D., Stevens, D.E.: The Robustness of honesty effects on budget proposals when the superior has rejection authority. Account. Rev. **90**(2), 467–493 (2015)
17. Blay, A., Douthit, J., Fulmer, B.: Why don't people lie? Negative affect intensity and preferences for honesty in budgetary reporting. Manag. Account. Res. **42**, 56–65 (2019)
18. Abdel-Rahim, H.Y., Stevens, D.E.: Information system precision and honesty in managerial reporting: a re-examination of information asymmetry effects. Account. Organ. Soc. **64**, 31–43 (2018)
19. Bicchieri, C.: Norms in the Wild: How to Diagnose, Measure, and Change Social Norms. Oxford University Press, New York, NY (2017)
20. Boster, C., Majerczyk, M., Tian, Y.: The effect of individual and pooled profit-sharing plans on honesty in managerial reporting. Contemp. Account. Res. **35**(2), 696–715 (2018)
21. Masclet, D., Noussair, C., Tucker, S., Villeval, M.-C.: Monetary and nonmonetary punishment in the voluntary contributions mechanism. Am. Econ. Rev. **93**(1), 366–380 (2003)
22. Noussair, C., Tucker, S.: Combining monetary and social sanctions to promote cooperation. Econ. Inq. **43**(3), 649–660 (2005)
23. Fehr, E., Gächter, S.: Cooperation and punishment in public goods experiments. Am. Econ. Rev. **90**(4), 980–994 (2000)
24. Evans, J.H., Hannan, R.L., Krishnan, R., Moser, D.V.: Honesty in managerial reporting. Account. Rev. **76**(4), 537–559 (2001)
25. Fehr, E., Fischbacher, U.: Third-party punishment and social norms. Evol. Hum. Behav. **25**(2), 63–87 (2004)
26. Dimant, E., Bicchieri, C., Xiao, E.: Deviant or wrong? The effects of norm information on the efficacy of punishment. SSRN (2018)
27. Cialdini, R.B., Reno, R.R., Kallgren, C.A.: A focus theory of normative conduct: recycling the concept of norms to reduce littering in public places. J. Pers. Soc. Psychol. **58**(6), 1015–1026 (1990)
28. Baiman, S.: Agency research in managerial accounting: a second look. Acc. Organ. Soc. **15**(4), 341–371 (1990)
29. Sutton, R.S., Barto, A.G.: Reinforcement learning. J. Cogn. Neurosci. **11**(1), 126–134 (1999)

# Modelling Regional Innovation Systems in Low and Middle-Income Countries

**Özge Dilaver, Alain Nkusi, Joshua Omoju, and Roseline Wanjiru**

**Abstract** Innovation is an important means of economic growth. It is, thus, much needed in low and middle income countries. Although this need was recognised early on, innovation literature is largely based on empirical studies in high-income countries. Theorising based on evidence on high-income countries can provide some insights about factors relevant to low and middle income countries if we assume countries follow similar growth trajectories. Yet, both innovation and growth are much more complex. Innovation occurs through complex interactions between various actors and it is often embedded in social contexts that vary across regions. It is, then, adopted and adapted to local needs and contexts. Agent-based social simulation offers an important potential to capture these complex dynamics. In this chapter, we put together two much-needed reviews: existing literature on innovation systems in low and middle income countries, and existing agent-based models of innovation systems. By juxtaposing our findings in these two reviews, we explore how agent-based models can address some of the major limitations of the information systems approach.

**Keywords** Innovation system · Regional innovation system · Low-income countries · Middle-income countries

## 1 Introduction: Innovation, Economic Growth and Innovation Systems

Innovation is thought to explain the observed differences between the growth of factors of production and the output. It is, therefore, considered as "the single, most important component of long-term economic growth" [58, p1]. Through the lens

Ö. Dilaver (✉) · J. Omoju · R. Wanjiru
Northumbria University, Newcastle Upon Tyne, UK
e-mail: ozge.dilaver@northumbria.ac.uk

A. Nkusi
Teesside University, Middlesbrough, UK

of methodological individualism of mainstream economics, innovation appears as the distinct outcome of firms' R&D investments. This level of abstraction, however, hides the inherent complexity of the innovation process that involves uncertainties and bounded rationality [19, 61]. Firms with similar assets and resources, operating in similar environments can hold different capabilities [46]. Thus, knowledge creation is path dependent and highly idiosyncratic [34, 49, 51]. Innovation systems (IS) perspective holds a more distributed nature of knowledge and events, and positions innovation as an interactive process between users and producers [43, 44], multiplicity of firms, as well as broader formal and informal institutional settings. Hence, it draws more complex and socially embedded representations of both innovation and economic growth that align with the methodological strengths of agent-based social simulation (ABSS).

Social embeddedness of innovation already implies the significance of places and geographies. Regions—areas that are distinguished from the surrounding land—captured interest in innovation studies because innovations are not uniformly distributed in space [8], and some regions significantly outperform others. Regional innovation systems (RIS) have long captured the interest of policy makers, who aim to emulate their success, and scholars, who aim to understand why businesses in related value chains agglomerate in certain places. The literature on specialised regions and industrial clusters (IC), dates further back at least to Marshall [45], who identified long-term benefits of knowledge exchange and diffusion of innovations between businesses and people when an industry "chosen a locality for itself" (p 27). RIS is a relatively new approach to innovation policies that view innovations as embedded within and shaped by dynamic interactions between actors in innovation networks, including private and public actors [6, 15]. It sees innovation as an interactive and cumulative process of learning and considers the critical role that strong regional capacities can contribute towards boosting the performance of systems of innovation [16, 29]. While some of the influential contributions to the study of RIS emerged from the work on the conceptual limitations and policy relevance of national innovation systems [8, 15, 16], and conceptual nuances can be identified between RIS and IC, they are similar enough to be used interchangeably (see, for example, [53]), and difficult to distinguish empirically. From a Schumpterian perspective, furthermore, since innovation is seen as endogenous to economic activity, it is difficult to separate production from innovation, and IC from RIS.

From this vantage point, existing literature on regions and innovativeness is voluminous, and at the core of this rich body of literature are case studies of successful regions in high-income countries. These studies highlight inter-firm exchanges and relationships. Firms are commonly thought as the main agents of innovation and firm-specific competencies and organisational learning are expected to create competitive advantages at the regional level due to their impact on skills, institutional environment and common assumptions and values.

Both IC and RIS are concerned with how geographical proximity facilitates economic growth and innovativeness. Although it was innovations in information and communication technologies that inspired many to imagine a placeless, border-less global economy, like many economic fields, innovation itself continues to

reveal spatial patterns. Proximity matters for knowledge production, and innovative processes, therefore, are often geographically concentrated [8, 68]. In this regard, it has been argued that innovation may be better understood and governed locally at regional level [7, 14].

While RIS of global importance and scale such as Silicon Valley, Detroit and Hollywood are investigated from various angles, and industrial clusters in OECD countries in general are observed via longitudinal analyses, our knowledge of innovation systems in low income countries remains scant. We do not know if innovation systems in low income countries impact upon global trade and global innovation systems. We also do not know how firms operating in industrial clusters in low income countries develop their dynamic capabilities, identify opportunities, respond to consumer affordability constraints and exchange knowledge across value chains.

If innovation is the main engine of economic development in market economies, these gaps in the existing literature are important for a broad range of stakeholders. In both low-income and middle-income countries, furthermore, effectiveness of innovation policies (at national and regional level) and strategies (at firm level) are of crucial importance due to limited resources. In this chapter, our main research question is: how can ABSS be useful for studying RIS in low-and middle-income (LMIC) countries? We address this question by taking stock of our existing understanding of innovation systems with a specific focus on RIS in low- and medium-income countries. We, then, review agent-based social simulation models of innovation models and explore research avenues for future models.

## 2 Conceptual and Methodological Limitations of Innovation Systems Approach

The innovation system approach improves our understanding of innovation as a complex, socially embedded and path-dependent process. Yet, it has some conceptual ambiguities and methodological limitations that may affect wider implementation of research findings. All of the terms used namely "innovation", "system", "national" [50] and "regional" [16] have been scrutinised. While the resources and activities of the firms that are at the world's technological frontiers can be different from others, strict definitions of innovation that are limited to these contexts can isolate innovations from their impact on economic growth, and so may not be always useful for research. Hence, innovation is often defined from the perspective of individual agents, whether or not something is new to the firm or the person is the criterion and not whether it is new to the world [50].

In a similar vein, "system" is an ambiguous concept. Systems can be naturally existing (e.g. ecosystem, cardiovascular system) or purposefully designed (e.g. information system, social welfare system). From an evolutionary point of view, furthermore, this important distinction may blur. Through mutation and survival both biological and societal systems can acquire function without any teleological design. Nelson

and Rosenberg [50] raise the issue of conscious design, and clarify that by "system", they refer to "*a set of institutional actors that, together, plays the major role in influencing innovative performance*" (p4). From a slightly different perspective Asheim and Gertler [8] highlight the "systemic", rather than idiosyncratic nature of relationships and practices within an innovation system. According to the authors, in order to be systemic, relationships between actors need to involve some degree of interdependency. As such, the term "system" is used loosely to capture interdependent, or game-theoretic relationships.

When it comes to "national", Nelson and Rosenberg [50] point out that on one hand this unit of analysis may seem too broad. On the other hand, it is too narrow and excludes the impact of globalisation. Elsewhere in social sciences, "nation" is already a contested term (see also [16]). Anderson [3] refers to "nation" as "an imagined community", a political invention of the nineteenth century. From this perspective, national territories are not well-defined containers of distinct economies and national borders often arbitrarily cut through social reality. On the other hand, states have unique powers in policy making, shaping educational institutions and distributing resources. "National" in this sense, is useful in capturing states' top-down impact on economic and innovative activities.

Despite being a core concept in multiple social science disciplines, "region" has multiple and ambiguous meanings too. The term is commonly used for defining areas at both sub-state (e.g. Italian districts) and supra-state levels. Region is, furthermore, 'ontologically slippery' [52], p. 21 in that it is not always clear if the region of interest exists independently of the attempts of studying it, or making sense of areas on a map. All of these conceptual ambiguities have implications for regional innovation systems, and in particular the bottom-up and top-down dynamics that are briefly mentioned above. Cooke et al. [16], therefore, distinguish between cultural and administrative regions, and relatedly processes of bottom up regionalism and top-down regionalisation.

Some of the methodological limitations of innovation systems approach stem from these conceptual difficulties. As the conceptual elements and properties of these systems are not developed a priori, research on RIS focuses on empirical descriptions of successful cases. The *ex poste* identification of structural, institutional, and organisational factors that are thought to have led to the emergence of more innovative, or competitive nations or regions falls short of determining the direction of causality, whether such factors impact success or vice versa. In addition, the broad scope of RIS and the rich contexts covered in case studies make it difficult to apply to other regions with different but equally rich contexts. Edquist [24] argues that innovation systems approach has limitations in providing formal propositions related to causality. As such it is more of a conceptual framework than a theory.

These conceptual and methodological shortcomings are particularly important for investigating RIS in LMIC. There are inherent challenges related to identification of innovation systems where the link between new ideas, technological applications and economic growth is much less visible. At the regional level furthermore, while it

is easy to identify and possible to empirically study complex and dynamic processes and networks of innovation in Silicon Valley, it is hard to know what to look for when such clear instances of RIS do not emerge.

These limitations on one hand, and the need for understanding RIS in LMIC on the other yield a peculiar form of knowledge accumulation. While the study of RIS in high-income countries focus on empirical case studies, arguably with a relative shortage of theoretical contributions, low and middle-income contexts are theorised without much empirical content. Authors theorise on innovation systems failures in comparison to ideal type innovation systems. In the context of RIS, for example, Tödtling and Trippl [66] argue that innovation activities differ strongly between central (metropolitan), peripheral and old industrial areas. Peripheral regions suffer from 'organisational thinness', they have less R&D intensity, incremental and process innovations instead of product innovations (see also [25, 57]).

## 3    What do we Know About Innovation in Low- and Middle-Income Countries?

This section aims to review and summarise our existing knowledge on RIS in low income countries. To our knowledge, there are no empirical studies of RIS in low income countries. We, therefore, expanded our scope to findings and arguments related to innovation and (both national and regional) innovation systems in LMIC. In doing so, we do not assume that findings on NIS and middle-income countries can be extended to RIS and low-income countries respectively, but we aim to take stock of what we know in these broader contexts to better identify remaining questions.

In studies on innovation in low income countries, the link to development has been widely acknowledged at least since the Sussex Manifesto [62] that called for strengthening scientific and technological capabilities of developing countries instead of transferring technologies from the developed countries or establishing hierarchical division of labour in science. This fault line between innovating in a way that suits local needs, infrastructures, skills and factor endowments, and acquiring technologies developed elsewhere remains important in more recent literature.

There is a growing number of studies that apply the systems perspective to low-income countries at the national level [12, 26, 36, 37, 39, 70]. One of the common arguments emerging from these studies is that the knowledge on NIS as conceived in high-income countries does not necessarily capture innovation processes and factors in low-income countries [37, 39] and so modifications are needed [26]. For example, given that innovation activities differ vastly between developed and developing economies [54], the effectiveness and relevance of indicators for measuring innovation are criticised for not putting adequate emphasis on the learning and capacity building processes [64], also have a look at: introducing innovation frameworks that transcends constraints imposed by contextual variations [39, 64]. Casadella and Tahi

[12] propose using qualitative indicators to capture learning capacity rather than focusing on R&D. Similarly, Khan [37] calls for considering inclusive absorptive capacity.

In a similar vein, regarding learning and developing capabilities, studies looking at low-income countries hold a common theme of *firms and catching up* [39]. They investigate innovation strategies of domestic firms in trying to upgrade their capabilities. However, since low-income countries experience critical barriers such as limited public resources to support innovation, challenges related to implementing policy in socio-economic conditions shaped by poverty, and relatedly, weak and/or short-term commitment of institutional actors [13], this approach often reports failures than success stories [4, 20, 66]. Relatedly, Kaplinsky and Kraemer-Mbula [36] identify two structural problems facing low and middle-income countries that may hinder their chances of growth through trajectories used in the past. These problems are the prevalence of the informal sector and the erosion of the growth by export model due to China's dominating role in international trade. The authors point out to new and multiple trajectories stemming from the innovative potential within the informal sector, the potential of South-South trade, and the transformative potential of new information and communication technologies in LMIC.

Lately, several authors argued the exclusion of the informal sector blurs our understanding of innovation systems in LMIC [26, 28]. For example, Egbetokun et al. [26] argues that innovation networks may involve more informal links than formalised ones in developing countries. Hence, the common assumption that regional actors in developing countries are poorly networked (see [7] may be overlooking informal links. In addition, in their study of the informal sector in Tanzania, Cozzens and Sutz's [18] found that the informal sector helps in adaptation of products produced elsewhere to local conditions in Tanzania. Furthermore, Fu et al. [28], who analysed patterns of innovation contributions in developing countries, found that although technological innovations have a greater overall contribution to firm productivity and growth than non-technological innovations, informal firms get as much out of non-technological innovations as formal firms do. The authors argue, therefore, that existing approaches to IS do not capture differences in the informal sector and the peculiar ways innovation occurs in such contexts.

As mentioned above, in our review we did not find any empirical studies on RIS in low-income countries. Some recent studies in lower-middle income countries in the Middle East and North Africa echo arguments in innovation literature. They point out that RIS in LMIC have features and dynamics that are qualitatively different from what is often referenced in the RIS literature [23, 27]. Djelfat and Cummings [23] highlight that even though supporting organizations are present, the critical linkages and the institutional set-ups that are needed to facilitate innovation, remain weak and fragmented. Hence, endogenous capabilities in these contexts are weak and well below what is required to address challenges in the productive sphere and society at large.

For these reasons, Djelflat and Cummings [23] propose that systems of innovation in developing countries and African countries in particular should be understood as systems in construction. This is similar to Fardj and Hammadi [27], who state that

the "emergence" paradigm is more appropriate to characterise innovation systems in neo-peripheral countries than the catch-up one. They cite evidence from Maghreb countries such as Morocco, where new proactive strategies for "emergence" are being undertaken, which includes the creation of zones and platforms dedicated to industrial structures, with an attractive incentive framework.

## 4  Agent-Based Social Simulation Models of Innovation Systems—An Overview

In this section, we provide an overview of the existing ABSS models of IS based on our systematic literature review. Our primary aim while conducting this review was better understanding existing RIS models, but we could only identify a few [11, 21, 30, 55, 57, 60, 65]. Using an approach similar to the previous section, we expanded our scope to gather insights from relevant models. Since ABSS is grounded in complexity theory and social networks have been modelled extensively through ABSS, some studies focused on innovation networks. These studies are relevant to our main research question in this chapter and so we included "networks" in our keywords. We scanned the existing literature using Web of Science and different combinations of keywords ("innovation systems", "innovation networks", or "innovation" and "agent-based modelling", or "agent-based modeling"). This scan produced our initial list of papers. As a secondary level of selection, we checked the papers that cited the 3 most highly cited papers in our initial list, and identified relevant studies.

In total, we identified 47 studies published as journal articles (39) or conference proceedings (8). These papers were published in 32 different publications, none of which published more than 2 of the papers in our list. In our view, this diversity indicates that there isn't yet an established dialogue, or a research area in ABSS models of IS. There is, however, a variety of approaches developing independently, and a growing number of studies (see, for example, [21, 38, 40] around the SKIN model [31, 32]).

In order to identify the similarities and differences between selected studies, we used the conceptual anatomy framework (CAF) proposed by Dilaver and Gilbert [22]. The framework represents five main conceptual elements of ABSS models namely, agents, environment, social structure, actions and interactions, and temporality. Our review both sheds light on the shared, or common aspects in ABSS models of IS, and demonstrates the diversity of these models in detail. In this respect, we expect this overview to contribute to the development of this research area, and the currently lacking dialogue between studies.

Within the five elements of the conceptual anatomy framework, the first element—agents—is most developed in the reviewed ABSS/IS studies in terms of richness and detail provided. Regarding types of agents, the typical agent in these models represents firms. To be clear, except Sebestyén and Varga [60] in which agents represent

regions, all of the studies we reviewed involved firm agents. [5, 10, 11, 30, 33, 69]. This is followed by universities [41, 47, 63, 65, 67]. Other types of agents include governments [55, 63], funding agencies [41], research agencies/centres ([65, 67], and venture capitalists [67]. In Dilaver et al. [21], in addition to organisations, there are people (researchers and managers). These agents exist in a multi-level structure: firms are an evolutionary output of the people they employ and those whom they employed in the past (see also [5]).

The heterogeneity of the firm agents is often represented via differences in knowledge. The "kene" concept from the SKIN model is influential in this conceptualisation. The kene represents the aggregate knowledge of a firm [1, 57, 67]. It consists of triples representing a firm's capability (C) in a scientific, technological, or business domain,its ability (A) to implement knowledge in this field; and its expertise level (E) with respect to the ability. In a similar vein, Dilaver et al. [21] place firms in a three-dimensional knowledge space. The horizontal dimensions of this space correspond to scientific areas and technological fields, and the vertical dimension corresponds to their level of expertise. Similarly, firms are represented as three-component vectors in Cannavacciuolo et al. [11] and this closely resembles the kene. Mao et al. [47] define four knowledge fields in which firms' knowledge levels vary. Slightly different from kene, the authors distinguish between explicit and implicit knowledge.

In a similar, but more dynamic way, heterogeneity is also thought in terms of organisational learning or absorptive capacity [11, 17, 47, 48, 65]. Relatedly, knowledge and capabilities are also approached in a relative manner, often expressed in terms of agents' cognitive distance from each other [9, 42] or network position [35]. Proximity, in this context, often corresponds to the probability and willingness for agents to form links [11, 60].

The model descriptions we reviewed were not as detailed and explicit for the second and third elements of the CAF framework; environment and social structure respectively. In most studies we reviewed, properties of the simulated environment and structures are described in-between lines that focus more on agents, or actions and interactions. Ponsiglione et al. [55], for example, introduce their competitive environment that holds regularities and relatedly innovation opportunities as an agent. In some studies, theoretical background and the model structure is explained together. In Cannavacciuolo et al. [11], the virtual environment is a closed system and this property is justified in reference to industrial districts literature. Some RIS models are based or calibrated with data on specific geographic areas [5, 9, 11, 60, 63, 65]. Dilaver et al. [21] study the effects of changing entrepreneurial character of regions on the development of industrial clusters in two simultaneously simulated regions based on the historical competition between Boston and Silicon Valley. Through 3D visualisation, the model accommodates three connected environments. One of these environments (earth space) represents local spatial dynamics, and the remaining two (reward space and knowledge space) represents innovation opportunities and knowledge resources of people and firms. Reward space furthermore is a semi-visible environment that each person agent sees differently according to their knowledge in related areas.

In terms of social structure, social networks are the most apparent feature of the studies we reviewed. Social networks are used both as explanatory [9, 47] and dependent variables [11, 35, 59]. Hence, on one hand, a priori network structure is thought to impact innovativeness, and on the other, interactions between agents are thought to yield *ex post* network properties. These dynamics indicate the evolutionary and path-dependent nature of relationships in IS. Kwon and Motohashi [38], furthermore, distinguishes between short-term and informal and long term and formal relationships.

In addition to networks, since most models we reviewed use firm agents, there are economic structures in the models we reviewed. The most common economic structure is a supply or value chain [11, 17, 30]. In some models, this structure is called a market [56, 57], yet it is not clear if these models entail demand, supply and price dynamics. In Dilaver et al.'s [21] multi-level structure, people agents are tied to firm agents through employment (as researchers or managers), or entrepreneurial actions (creating a start-up/spinoff).

Regarding actions and interactions, the fourth CAF element, the most important dynamic is the choice of a partner to collaborate, or to exchange (buy or sell products). These dynamics reflect the assumption of bounded rationality in that agents do not fully know the variety of capability and the knowledge level of other firms [30, 60, 67]. The models we reviewed compared the impact of different search and match algorithms based on random selection, spatial selection, selection based on knowledge capital and/or complementarity of knowledge [2, 40, 57].

Finally, in terms of temporality, the fifth element of the CAF framework, the model descriptions we reviewed are not very detailed. Time dimension is not explicitly discussed. In most of these models, innovations are key events that impact upon simulated histories. Innovation is often (see, for example, [21, 35] thought of as a function of various knowledge requirements and some models describe the rules that determine innovation events in detail and with graphics. Through above mentioned interactions, firms acquire and accumulate knowledge over time and when this accumulation satisfies pre-set conditions, successful innovations emerge. In this context, innovation represent different outputs including scientific papers, patents [65, 67], Mahmoudzadeh an Alborzi, 2017, quality of products [30], profits [57], or the number of firms as an indirect measure of profits [21] and knowledge creation [47] or diffusion [9, 40]. As these are related but different concepts, explicit discussions of simulated histories can support cross-fertilisation between ABSS models of IS.

## 5 Conclusions

RIS perspective regards innovation as an interactive process that involves both bottom-up and socially embedded dynamics and the impact of top-down policies. Compared to methodological individualism of mainstream economics, this way of approaching innovation and economic growth aligns better with the methodological

strengths of ABSS. Interaction within and between users and producers involve interdependencies and path, or history-dependency. Knowledge, furthermore, has a tacit, subjective and distributed nature that calls for special attention to heterogeneities and their evolution over time.

While the need for understanding RIS in contexts where economic growth is most needed was identified early on, due to conceptual and methodological limitations constraining empirical research, the literature on IS in general and on RIS in particular is generally limited to repeated calls for further research on one hand, and theoretical analyses of "system failure" developed through comparisons to ideal type RIS. Overall, the established indicators and measures that are largely driven from high-income countries fall short of capturing innovation and growth in LMIC where regional transformations would be most needed. Our review of the existing literature highlighted three directions for producing research relevant in this context. Firstly, echoing the Sussex Manifesto (see Sect. 2), there is a need for studying innovation in low and middle-income countries in their own context, in terms of developing capabilities for addressing local needs, or needs in similar LMIC. Secondly and relatedly, there is a need for moving beyond tautological explanations of RIS "failure" (e.g. poorer regions are poor because they are poor). Finally, the role of informal, or shadow economy in learning, innovation and growth requires fresh and unbiased perspectives.

Existing ABSS models of IS offer several strengths for addressing these needs in *silico.* These models can represent learning and innovation at multiple levels in a way that supports the broad conception of innovation at the level of individuals and organisations, while also making it possible to follow the links between micro and macro events. They can generate stochastic, path-dependent processes that better represent multiplicity of growth trajectories compared to overly deterministic catching-up models. ABSS/IS models can also support improved understandings of self-reproducing patterns of learning, networks and capital. Finally, these models can be used for capturing the semi-visible and fuzzy dynamics of informal economies.

There are also some challenges for ABSS models in capturing RIS in LMIC. Although the relevant literature is growing, this growth is occurring in silos. More research effort is needed to increase awareness of existing studies, enable cross-fertilisation and create a dialogue between these studies. We believe our review (Sect. 4) will contribute towards improving the transparency of ABSS IS models. A second challenge is that like the IS field in general, some of the influential ABSS models are built with IS in high-income countries in mind. Hence, in modelling of RIS in LMIC, the relevance of existing measures and representations need to be carefully reconsidered in the light of relevant literature that we reviewed (Sect. 3) in this chapter.

# References

1. Ahrweiler, P., Gilbert, N., Pyka, A.: Agency and structure: a social simulation of knowledge-intensive industries. Comput. Math. Organ. Theory **17**(1), 59–76 (2011)
2. Akbas, M.I., Gunaratne, C., Garibay, O.O., Garibay, I., O'Neal, T.: Role of entrepreneurial support for networking in innovation ecosystems: an agent based approach. In: 2015 Winter Simulation Conference (WSC) (2015)
3. Anderson, B.: Imagined Communities: Reflections on the Origin and Spread of Nationalism. Verso, London (1991[1983])
4. Arocena, R., Sutz, J.: Weak knowledge demand in the South: learning divides and innovation policies. Sci. Pub. Policy **37**(8), 571–582 (2010)
5. Arutyunovich, A.S., Yuryevich, A.M.: Assessment of innovative activity of regions in the Russian Federation. Montenegrin J. Econ. **11**(1), 7–21 (2015)
6. Asheim, B.T., Isaksen, A., Trippl, M.: Advanced Introduction to Regional Innovation Systems (2019)
7. Asheim, B., Coenen, L., Moodysson, J.: Methods and applications of regional innovation systems analysis. In: Handbook of Research Methods and Applications in Economic Geography, pp. 272–290. Edward Elgar Publishing (2015)
8. Asheim, B.T., Gertler, M.S.: The geography of innovation: regional innovation systems. In: Fagerberg, J., Mowery, D.C. (eds.) The Oxford Handbook of Innovation Oxford Academic, pp. 291–317 (2009). https://doi.org/10.1093/oxfordhb/9780199286805.003.0011.
9. Bogner, K., Müller, M., Schlaile, M.P.: Knowledge diffusion in formal networks: the roles of degree distribution and cognitive distance. Int. J. Comput. Econ. Econometrics **8**(3/4), 388 (2018)
10. Caiani, A.: Innovation dynamics and industry structure under different technological spaces. Ital. Econ. J. **3**(3), 307–341 (2017)
11. Cannavacciuolo, L., Iandoli, L., Ponsiglione, C., Zollo, G.: Learning by failure VS learning by habits. Int. J. Entrep. Behav. Res. **23**(3), 52 (2017)
12. Casadella, V., Tahi, S.: national innovation systems in low-income and middle-income countries: re-evaluation of indicators and lessons for a learning economy in Senegal. J. Knowl. Econ., pp. 1–31 (2022)
13. Chaminade, C., Padilla-Pérez, R.: The challenge of alignment and barriers for the design and implementation of science, technology and innovation policies for innovation systems in developing countries. In: Research Handbook on Innovation Governance for Emerging Economies, pp. 181–204. Edward Elgar Publishing (2017)
14. Cooke, P.: Regional innovation systems, clusters, and the knowledge economy. Ind. Corp. Chang. **10**(4), 945–974 (2001)
15. Cooke, P.: Regional innovation systems: competitive regulation in the new Europe. Geoforum **23**(3), 365–382 (1992)
16. Cooke, P., Uranga, M.G., Etxebarria, G.: Regional innovation systems: institutional and organisational dimensions. Res. Policy **26**(4–5), 475–491 (1997)
17. Cozzoni, E., Passavanti, C., Ponsiglione, C., Primario, S., Rippa, P.: Interorganizational collaboration in innovation networks: an agent based model for responsible research and innovation in additive manufacturing. Sustainability **13**(13), 7460 (2021)
18. Cozzens, S., Sutz, J.: Innovation in informal settings: reflections and proposals for a research agenda. Innov. Dev. **4**(1), 5–31 (2014)
19. Cyert, R.M., March, J.G.: A Behavioral Theory of the Firm, vol. 2, no. 4, pp. 169–187. Englewood Cliffs, NJ (1963)
20. Delvenne, P., Thoreau, F.: Dancing without listening to the music: learning from some failures of the 'national innovation systems' in Latin America. In: Research Handbook on Innovation Governance for Emerging Economies, pp. 37–58. Edward Elgar Publishing (2017)
21. Dilaver, Ö., Bleda, M., Uyarra, E.: Entrepreneurship and the emergence of industrial clusters. Complexity **19**(6), 14–29 (2014)

22. Dilaver, Ö., Nigel, G.: Unpacking a black box: a conceptual anatomy framework for agent-based social simulation models. J. Artif. Soc. Soc. Simul. **26**(1), 4. http://jasss.soc.surrey.ac.uk/26/1/4.html. https://doi.org/10.18564/jasss.4998(2023)
23. Djeflat, A., Cummings, A.: Emergence of territorial systems of innovation in developing countries: building a conceptual framework based on Latin American and North African experiences (2015)
24. Edquist, C.: Systems of innovation: perspectives and challenges. In Fagerberg, J., Mowery, D.C. (eds.) The Oxford Handbook of Innovation Oxford Academic (2009)
25. Edquist, C., Eriksson, M.L., Sjögren, H.: Characteristics of collaboration in product innovation in the regional system of innovation of East Gothia. Eur. Plan. Stud. **10**(5), 563–581 (2002)
26. Egbetokun, A., Oluwadare, A.J., Ajao, B.F., Jegede, O.O.: Innovation systems research: an agenda for developing countries. J. Open Innov. Technol. Market Complexity **3**(4), 1–16 (2017)
27. Ferdj, Y., Hammadi, A.: Dynamic and emergence of development territorial, Algerian cluster study. Revue tadamsa d'unegmu **3**(1), 96–108 (2023)
28. Fu, X., Mohnen, P., Zanello, G.: Innovation and productivity in formal and informal firms in Ghana. Technol. Forecast. Soc. Change **131**(C), 315–325 (2018)
29. Fernandes, C., Farinha, L., Ferreira, J.J., Asheim, B., Rutten, R.: Regional innovation systems: what can we learn from 25 years of scientific achievements? Reg. Stud. **55**(3), 377–389 (2021)
30. Giardini, F., Di Tosto, G., Conte, R.: A model for simulating reputation dynamics in industrial districts. Simul. Model. Pract. Theory **16**(2), 231–241 (2008)
31. Gilbert, N., Pyka, A., Ahrweiler, P.: Innovation networks—a simulation approach. J. Artif. Soc. Soc. Simul. **4**(3). http://www.soc.surrey.ac.uk/JASSS/4/3/8.html (2001)
32. Gilbert, N., Ahrweiler, P., Pyka, A.: Learning in innovation networks: some simulation experiments. Physica A **378**(1), 100–109 (2007)
33. Heshmati, A., Lenz-Cesar, F.: Policy simulation of firms' cooperation in innovation. Res. Eval. **24**(3), 293–311 (2015)
34. Hodgson, G.M.: Competence and contract in the theory of the firm. J. Econ. Behav. Organ. **35**(2), 179–201 (1998)
35. Hua, L., Wang, W.: The impact of network structure on innovation efficiency: an agent-based study in the context of innovation networks. Complexity **21**(2), 111–122 (2015)
36. Kaplinsky, R., Kraemer-Mbula, E.: Innovation and uneven development: the challenge for low-and middle-income economies. Res. Policy **51**(2), 104394 (2022)
37. Khan, M.S.: Absorptive capacities approaches for investigating national innovation system in low and middle income countries. Int. J. Innov. Stud. **6**(3), 183–195 (2022)
38. Kwon, S., Motohashi, K.: How institutional arrangements in the national innovation system affect industrial competitiveness: a study of Japan and the US with multiagent simulation. Technol. Forecast. Soc. Change **115**, 221–235 (2017)
39. Lema, R., Kraemer-Mbula, E., Rakas, M.: Innovation in developing countries: examining two decades of research. Innov. Dev. **11**(2–3), 189–210 (2021)
40. Li, L., Xie, J., Wang, R., Su, J., Sindakis, S.: The partner selection modes for knowledge-based innovation networks: a multiagent simulation. IEEE Access **7**, 140969–140979 (2019)
41. London, J.O.N., Sheikh, N.J.: Innovation in African-American high-tech enterprises: a multi-agent approach. Entrepreneurship Sustain. Issues **7**(4), 3101 (2020)
42. Long, Q., Li, S.: A multi-agent-based evolution model of innovation networks in dynamic environments. In: 2014 International Conference on Mathematics and Computers in Sciences and in Industry (2014). https://doi.org/10.1109/mcsi.2014.34
43. Lundvall, B. (ed.): National Systems of Innovation: Towards a Theory of Innovation and Interactive Learning. Pinter, London (1992)
44. Lundvall, B.: Innovation as an interactive process: from user-producer interaction to the national system of innovation. In: Dosi, G., et al. (eds.) Technical Change and Economic Theory, pp. 349–369. Pinter, London (1988)
45. Marshall, A.: Principles of Economics: An Introductory Volume. Macmillan, London (1920)
46. Malerba, F., Orsenigo, L.: Knowledge, innovative activities and industrial evolution. Ind. Corp. Change. **9**(2), 289–314 (2000)

47. Mao, C., Yu, X., Zhou, Q., Harms, R., Fang, G.: Knowledge growth in university-industry innovation networks—results from a simulation study. Technol. Forecast. Soc. Change **151** (2020)
48. Muller, M., Kudic, M., Vermeulen, B.: The influence of the structure of technological knowledge on inter-firm R&D collaboration and knowledge discovery: an agent-based simulation approach. J. Bus. Res. **129**, 570–579 (2021)
49. Nelson, R.R.: Why do firms differ, and how does it matter? Strateg. Manag. J. **12**(S2), 61–74 (1991)
50. Nelson, R.R., Rosenberg, N.: Technical innovation and national systems. In Nelson, R.R. (ed.) National Innovation Systems: A Comparative Analysis. Oxford University Press, New York, pp. 3–22 (1993)
51. Nonaka, I., Toyama, R., Nagata, A.: A firm as a knowledge-creating entity: a new perspective on the theory of the firm. Ind. Corp. Chang. **9**(1), 1–20 (2000)
52. Paasi, A., Metzger, J.: Foregrounding the region. Reg. Stud. **51**(1), 19–30 (2017)
53. Padmore, T., Gibson, H.: Modelling systems of innovation: II. A framework for industrial cluster analysis in regions. Res. Policy **26**(6), 625–641 (1998)
54. Pfotenhauer, S.M., Wentland, A., Ruge, L.: Understanding regional innovation cultures: narratives, directionality, and conservative innovation in Bavaria. Res. Policy **52**(3), 104704 (2023)
55. Ponsiglione, C., Quinto, I., Zollo, G.: Regional innovation systems as complex adaptive systems: the case of lagging European regions. Sustainability **10**(8), 2862 (2018)
56. Pyka, A., Gilbert, N., Ahrweiler, P.: Simulating knowledge-generation and distribution processes in innovation collaborations and Networks. Cybern. Syst. **38**(7), 667–693 (2007)
57. Pyka, A., Kudic, M., Müller, M.: Systemic interventions in regional innovation systems: entrepreneurship, knowledge accumulation and regional innovation. Reg. Stud. **53**(9), 1321–1332 (2019)
58. Rosenberg, N.: Innovation and Economic Growth in Innovation and Growth in Tourism, OECD Publishing, Paris. https://doi.org/10.1787/9789264025028-4-en (2004)
59. Savin, I., Egbetokun, A.: Emergence of innovation networks from R&D cooperation with endogenous absorptive capacity. J. Econ. Dyn. Control **64**, 82–103 (2016)
60. Sebestyén, T., Varga, A.: Knowledge networks in regional development: an agent-based model and its application. Reg. Stud. **53**(9), 1333–1343 (2019)
61. Simon, H.A.: From substantive to procedural rationality. In: 25 Years of Economic Theory: Retrospect and Prospect, pp. 65–86 (1976)
62. Singer, H., Cooper, C., Desai, R.C., Freeman, C., Gish, O., Hill, S., Oldham, G.: Draft introductory statement for the world plan of action for the application of science and technology to development, Annex II in Science and Technology for Development: Proposals for the Second Development Decade, United Nations, Dept of Economic and Social Affairs, New York, Document ST/ECA/133 (1970)
63. Sjafrina, N., Udin, F., Anggraeni, E.: A mapping of current downstream shallot supply chain based on agent-based modeling and quadruple innovation helix: a case study at Cirebon district, Indonesia. In: IOP Conference Series: Earth and Environmental Science, vol. 472, no. 1, p. 012056, Apr 2020. IOP Publishing (2020)
64. Sutz, J.: Measuring innovation in developing countries: some suggestions to achieve more accurate and useful indicators. Int. J. Technol. Learn. Innov. Dev. **5**(1–2), 40–57 (2012)
65. Uribe-Gómez, J.A., Giraldo-Ramírez, D.P., Gallón-Londoño, L., Fernandez-Ledesma, J.D., Quintero-Ramírez, S.: Analysis of dynamics, structures and agent relationships in regional innovation systems. Estudios Gerenciales **35**(153), 379–387 (2019)
66. Todtling, F., Trippl, M.: One size fits all? Towards a differentiated regional innovation policy approach. Res. Policy **34**(8), 1203–1219 (2005)
67. Triulzi, G., Pyka, A.: Learning-by-modeling: insights from an agent-based model of university-industry relationships. Cybern. Syst. **42**(7), 484–501 (2011)
68. Uyarra, E.: What is evolutionary about 'regional systems of innovation'? Implications for regional policy. J. Evol. Econ. **20**(1), 115 (2010)

69. van Rijnsoever, F.J.: Meeting, mating, and intermediating: how incubators can overcome weak network problems in entrepreneurial ecosystems. Res. Policy **49**(1), 103884 (2020)
70. Wandera, F.H.: The innovation system for diffusion of small wind in Kenya: strong, weak or absent? A technological innovation system analysis. Afr. J. Sci. Technol. Innov. Dev. **13**(5), 527–539 (2021)

# Embedding Social Simulation in the Design of Wine Pricing Policies

Nikitas M. Sgouros

**Abstract**  We provide an overview of Politika, a policy design prototype, and explain how it is applied in developing and analyzing pricing policies for wine brands versus their competitors. These policies seek to maximize the purchase motivation for specific brands of wine relative to their competitors in a population. Politika provides explicit representations for the policy parameters and their base case values that reflect the current state of the market. It then represents each policy alternative as a set of alternative values for a subset of the policy parameters. Furthermore, it is able to describe a set of constraints in the simulation of each alternative that can facilitate comparisons between alternatives. Finally, Politika allows the definition of criteria that will be automatically applied to the outcomes of the simulations and will allow the designer to estimate whether each alternative fulfills the policy goals.

**Keywords**  Policy design · Agent-based simulation · Consumer behavior

## 1  Introduction

One of the most frequent and difficult problems in policy design is that no empirical data exist for the outcomes of many of the proposed policy alternatives. For example, our knowledge of the effects of different pricing of particular wines in specific markets come from empirical data only for the pricing policies that have been applied to these products previously. In this case, simulation provides a promising solution to the analysis of such untried alternatives justifying the need for it to become tightly integrated to policy design. Furthermore, policy design can result in the formulation of incompatible policy alternatives. Thus, it is important to develop design and analysis tools that provide transparency and facilitate comparisons between policy alternatives. We provide an overview of Politika, a policy design prototype and its

N. M. Sgouros (✉)
Department of Digital Systems, University of Piraeus, Piraeus18534, Greece
e-mail: sgouros@unipi.gr

application to the design of wine pricing policies. These policies seek to increase the purchase motivation of selected wine brands in relation to their competitors in a specific population. Politika forms the social dynamics component of the Policy-CLOUD architecture [3], an EU-funded project that aims to provide analytic tools for supporting policy modeling and design on the cloud.

## 2   Modeling of Wine Purchase Motivation

Each simulation model in Politika includes:

1. A list of policy attributes and their initial values.
2. A set of rules for policy dynamics. These rules specify. how the policy attributes change during the simulation.
3. A list of individual attributes and their initial values. These describe the characteristics of each individual agent in the population.
4. A set of rules for individual dynamics. These describe how attributes for each individual are updated during the simulation.
5. A set of rules for connection dynamics. These describe how both connections and their attributes change during the simulation. Such dynamics depend on the attributes of the nodes at both ends of the connections and on the policy-relevant attributes.
6. The size of the population on which the policy will be simulated. along with the number of cycles for which the simulator will run.
7. The specifications of a network generator component. This is used whenever the user wants to simulate policy effects on a population generated from a known network model such as a random graph or various power-law networks. The features of such a model are provided as a set of key-value pairs.
8. A list of metric units used for all policy, individual or connection attributes.

In developing a simulation model for wine purchase motivation, we assume that price and quality are the main factors influencing consumers when purchasing wine. In addition consumers can be influenced by their exposure to wine-related advertising/marketing campaigns and the wine preferences of their social circle. Based on these assumptions we define the following set of parameters of interest for estimating the purchase motivation for a particular brand of wine (e.g., A) in a specific region:

1. Actual price for A
2. Quality (in a scale of 0 to 1) of A as determined by its average rating in a series of online reviews.
3. Estimate of the average price of wines sold in the region of interest.
4. Estimate of the maximum price for wine that is acceptable for an average consumer (e.g. double the average price of wines sold in the region).
5. Average quality of the wines sold in the region of interest (0 to 1).
6. Average income of the population in the region of interest.

7. Upper income of the population in the region of interest (e.g., double the value of the mean income in the population).
8. Average relative exposure of individuals to the advertising campaign for A (0 to 1). We assume that average exposure is proportional to the relative size of the advertising budget for A compared to its competitors.

We further assume that the population in the region of interest is represented as a social network, where each node corresponds to an individual. For each individual (e.g., X), each outgoing edge is labeled with a weight (`contact_strength`) representing the influence that X exerts on the wine purchasing decisions of one of its social connections. The purchase influence of X towards A (`purchase_influence_A`) is computed as the product of the `contact_strength` of the current edge times the current purchase motivation of X towards A (`purchase_motiv_A`).

Each individual X has a set of attributes that are relevant towards A. These include X's:

1. Income ranking (in a scale of 0 to 1) as determined by the ratio of its income to the maximum income for the region.
2. Sensitivity to the price of A as determined by the product of the difference of 1 minus X's income ranking times the ratio of the current price of X to the maximum wine price in the region. Therefore, price sensitivity provides an estimate of how much the price of A affects X's willingness to buy it. According to this estimate, poor individuals are more sensitive to the price of wines compared to wealthier ones.
3. Sensitivity to the quality of A as determined by the ratio of the current quality of A to the average quality of wines in the region, times X's income ranking. Therefore, quality sensitivity provides an estimate of how much the quality of A affects X's' willingness to buy it. According to this estimate, wealthier individuals are more sensitive to the quality of wines than poor ones.
4. Susceptibility to the advertising/marketing campaign for A (0 to 1). This estimates the extent to which an individual attends to and values ad messages as sources of information for guiding her consumptive behavior. This can depend on the exposure of X to the ad campaign with more exposure leading to less susceptibility.
5. Perceived influence for A from X's social circle. This is computed as the average purchase influence for A stemming from X's social circle.

Based on these attributes, the model estimates X's purchase motivation for A as a linear combination of:

1. X's price sensitivity for A.
2. X's quality sensitivity for A.
3. The product of X's advertising susceptibility for A to the intensity of A's ad campaign. We assume that the intensity of the ad campaign for A is a real number between 0 and 1 that is proportional to the relative size of the advertising budget for A compared to its competitors but also to the type of ad campaign (e.g. targeted or undirected) used.

4. The perceived influence for A from X's social circle.

Figure 1 describes the coding of the base case for the wine pricing policy model in Politika. This is the case that reflects the current status between two competing wines A and B in terms of price, quality and ad intensity for a specific population of 1000 individuals with an average number of three connections. The "1" specification in the units denotes unitless parameters.

## 3 Creation and Comparison of Wine Pricing Alternatives

Politika associates with a policy a set of design constraints that include:

1. **The base case scenario**. This is the set of policy attributes and their values that reflect the current state of the world before any policy intervention is tried.
2. **The set of policy alternatives that will be explored during design**. These are denoted as alternative sets of attribute-value pairs that denote the subset of policy attributes and their values that are different from the base case for each alternative.
3. **The population model relevant to the policy**. We assume that the population can be described as a graph in which individuals correspond to nodes and their relations as edges. In order to facilitate comparisons and reason about policy alternatives,it is often the case that their analysis should be based on the same population model. Consequently, the definition of the population model is adopted by all policy simulations in their network generator component when creating their population of individuals.
4. **The number of simulation rounds and sizes of populations on which each alternative will be simulated**.
5. **A set of criteria for evaluating the outcome of each alternative**. Typically such criteria are desired values for some policy attributes after the simulation of each alternative or desired ratios between such attributes or a combination of the two. The application of a criterion on the outcomes for an alternative will result in a true/false value.

Based on these constraints Politika automatically generates a bottom-up processing pipeline for transforming simulation outcomes of the various alternatives into policy recommendations. More specifically, when the user chooses to simulate all policy alternatives associated with a specific simulation model then Politika applies the design constraints defined for it and automatically runs all the simulations involved. It stores the results of all the rounds of simulations indexed under the round number for each. These results are then used to compute a set of analytics for each of the policy-relevant attributes defined in the design. Currently, this set includes the average value of each attribute along with its minimum and maximum value after all the simulation rounds for each alternative.

In order to design a policy for improving the purchase motivation of a specific wine versus a competitor in Politika, the policy maker provides data from specialist

- **Size:** 100

- **Generator:** max_edges: 5, min_edges: 1, directed?: false, degree_distribution: :homogeneous

- **Policy related parameters:** max_price: 30, price_A: 13, price_B: 11, avg_price: 15, quality_A: 0.87, quality_B: 0.89, avg_quality: 0.85, max_income: 50000, ad_intensity_A: 0.2, ad_intensity_B: 0.2, avg_purchase_motivation_A: 0, avg_purchase_motivation_B: 0, purchase_motivation_ratio: 0, wine_price_ratio: -1, ad_intensity_ratio: -1

- **Units for Policy related parameters:** max_price: "Euros", price_A: "Euros", price_B: "Euros", avg_price: "Euros", quality_A: "1", quality_B: "1", avg_quality: "1", max_income: "Euros", ad_intensity_A: "1", ad_intensity_B: "1", avg_purchase_motivation_A: "1", avg_purchase_motivation_B: "1", purchase_motivation_ratio: "1", wine_price_ratio: "1", ad_intensity_ratio: "1"

- **Individual attributes:** income_ranking: 0, purchase_motiv_A: 0, purchase_motiv_B: 0, price_sensitivity_A: -100, price_sensitivity_B: -100, quality_sensitivity_A: -100, quality_sensitivity_B: -100, ad_susceptibility: rand_float(0, 1), perceived_influence_A: 0, perceived_influence_B: 0, income: max(:rand.normal(30000, 5000*5000), 1)

- **Units for Individual attributes:** income_ranking: "1", purchase_motiv_A: "1", purchase_motiv_B: "1", price_sensitivity_A: "1", price_sensitivity_B: "1", quality_sensitivity_A: "1", quality_sensitivity_B: "1", ad_susceptibility: "1", perceived_influence_A: "1", perceived_influence_B: "1", income: "€"

- **Connection attributes:** contact_strength: rand_float(0, 1), purchase_influence_A: 0, purchase_influence_B: 0, red: rand_int(0, 255), green: rand_int(0, 255), blue: rand_int(0, 255)

- **Units for Connection attributes:** contact_strength: "1", purchase_influence_A: "1", purchase_influence_B: "1", red: "1", green: "1", blue: "1"

- **Policy related dynamics:** IF true DO this.avg_purchase_motivation_A: global_aggregate("this.purchase_motiv_A", population)/size; this.avg_purchase_motivation_B: global_aggregate("this.purchase_motiv_B", population)/size; this.purchase_motivation_ratio: global_aggregate("this.purchase_motiv_A", population)/global_aggregate("this.purchase_motiv_B", population) ~ IF this.ad_intensity_ratio == -1 and this.ad_intensity_B != 0 DO this.ad_intensity_ratio: this.ad_intensity_A/this.ad_intensity_B ~ IF this.wine_price_ratio == -1 and this.price_B != 0 DO this.wine_price_ratio: this.price_A/this.price_B

- **Individual dynamics:** IF this.income > global.max_income DO this.income: global.max_income ~ IF this.income <= global.max_income DO this.income_ranking: this.income/global.max_income ~ IF true DO this.perceived_influence_A: rand_float(0, 1); this.perceived_influence_B: rand_float(0, 1) ~ IF this.income_ranking != 0 and global.avg_price != 0 and global.price_A != 0 and global.price_B != 0 DO this.price_sensitivity_A: (global.price_A/global.max_price)*(1-this.income_ranking); this.price_sensitivity_B: (global.price_B/global.max_price)*(1-this.income_ranking) ~ IF this.income_ranking != 0 and global.avg_quality != 0 DO this.quality_sensitivity_A: (global.quality_A/global.avg_quality)*this.income_ranking; this.quality_sensitivity_B: (global.quality_B/global.avg_quality)*this.income_ranking ~ IF this.price_sensitivity_A != -100 and this.price_sensitivity_B != -100 and this.quality_sensitivity_A != -100 and this.quality_sensitivity_B != -100 DO this.purchase_motiv_A: -this.price_sensitivity_A + this.quality_sensitivity_A + this.ad_susceptibility * global.ad_intensity_A + sum_in(this_state, "edge.purchase_influence_A")/in_degree(this_state); this.purchase_motiv_B: -this.price_sensitivity_B + this.quality_sensitivity_B + this.ad_susceptibility * global.ad_intensity_B + sum_in(this_state, "edge.purchase_influence_B")/in_degree(this_state)

- **Connection dynamics:** IF this.purchase_motiv_A != 0 and this.purchase_motiv_B != 0 DO edge.purchase_influence_A: edge.contact_strength * this.purchase_motiv_A; edge.purchase_influence_B: edge.contact_strength * this.purchase_motiv_B;

**Fig. 1** Specification of the simulation model for the base case policy in our example for two competing wines A and B

**Compose Policy based on Ignacio Marín 2019 Magnifico Garnacha (Cariñena) VS Torres 1996 Gran Sangre de Toro Reserva Red (Penedès) - 3D**

**Title:**

Price Discovery for Aragon vs Catalan Wine

**Description:**

Find pricing and ad intensity ratios for an Aragon wine vs a competitor that can provide higher purchase motivation for the Aragon wine than its competitor. Assume that purchase motivation depends on price, quality ad intensity and social influence.

**Base Case:**

max_price: 30, price_A: 13.0, price_B: 11.0, avg_price: 15, quality_A: 0.87, quality_B: 0.89, avg_quality: 0.85, max_income: 50000, ad_intensity_A: 0.2, ad_intensity_B: 0.2, avg_purchase_motivation_A: 0, avg_purchase_motivation_B: 0, purchase_motivation_ratio: 0, wine price ratio: -1, ad intensity ratio: -1, x: 0, y: 0, z: 0

**Alternatives:**

price_A: 11, ad_intensity_A: 0.4 | price_A: 9, ad_intensity_A: 0.2

**Criteria:**

adequacy "avg_purchase_motivation_A.avg > avg_purchase_motivation_B.avg",  effectiveness: "(avg_purchase_motivation_A.avg/avg_purchase_motivation_B.avg) > 1.1"

**Social Network Generator:**

max_edges: 6, min_edges: 1, directed?: false, degree_distribution: :homogeneous
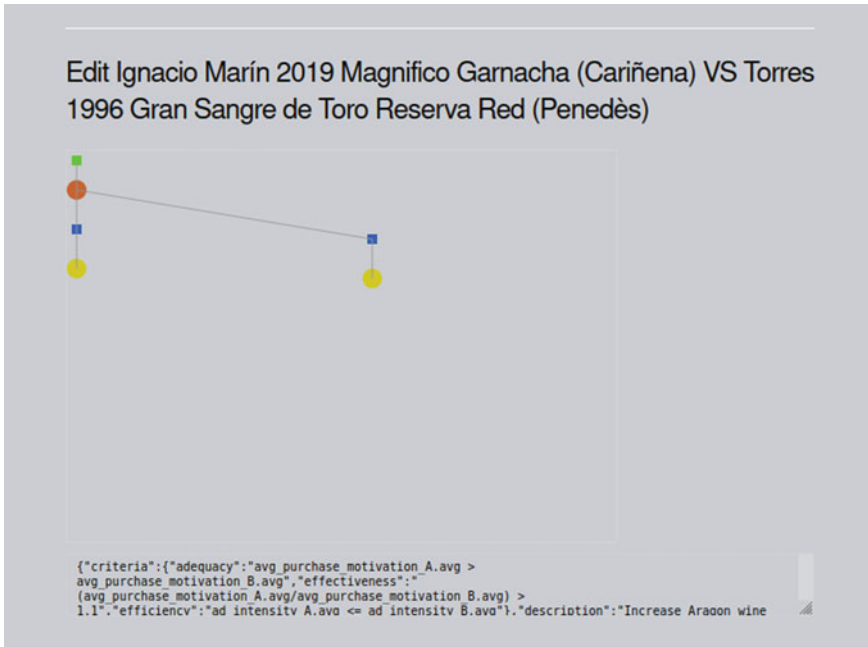
**Population:** 1000                                            **Repetitions:** 2

**Fig. 2** Specification of design constraints and alternatives for pricing policy (A = Ignacio Marin, B = Torres)

wine sites (e.g. prices, quality) and Wikipedia (e.g. average income for the population of interest) for two competing wines, e.g. A and B, for the parameters of interest in our simulation model.

The policy maker can then define and simulate various alternatives for pricing and/or advertisement effort for A in order to discover the mix that could improve $A's$ average purchase motivation in the population with respect to B in a specific population. For example,in Fig. 2 the user has defined two alternatives ($price_A$ : 11, $ad_intensity_A$ : 0.4 or $price_A$ : 9, $ad_intensity_A$ : 0.2) separated with $'|'$. In the same Figure the user has defined two criteria for evaluating each alternative. The adequacy criterion is satisfied when the average purchase motivation for A is found to be greater than the one for B, while the effectiveness criterion is satisfied if the ratio of the motivation for A versus B is greater than 1.1. Each simulation will run on

**Fig. 3** Tree-based GUI for the specification of design constraints and alternatives for our wine pricing policy

a social network generated using the specifications in the Social Network Generator field where all nodes will have a uniformly distributed degree between 1 and 6 and the edges will be undirected. The size of each network will be 1000 nodes and there will be 2 simulation rounds (repetitions) for each alternative.

We validate the computed relation ($>$, $=$, $<$) of the purchase motivation for A versus B in the base case using as proxy the relation between the number of ratings for them in such wine specialist sites (if A has a higher number of ratings than B this means that more people have purchased A therefore we can assume that currently the purchase motivation for A is greater than the one for B).

Politika provides an alternative interface for modeling policies through a tree-based GUI (see Fig. 3). The root of such a tree contains a Policy node having a set of alternative Goals as its children. Each Goal contains an abstract description of the desired outcomes of a policy. Under each Goal hangs a set of alternative Objectives for achieving this Goal. An Objective corresponds to a specific methodology for achieving a goal. It represents a policy alternative for a specific Goal. Each Objective, in turn, is decomposed into a sequence of Steps. Each Step represents a policy execution step in the methodology of the parent Objective. We assume that the execution of each Step can be simulated, thus providing a value range for its possible outcomes. Figure 3 visualizes the policy as a tree structure in which the green (root) node contains the description of the Policy and the red node contains the description

**Table 1** Results for the design of a pricing policy for wine X versus Y

| # | Price | Ad intensity | Purchase motivation |
|---|---|---|---|
| 1 | 0.81818 | 1 | 1.026 |
| 2 | 1.1 | 1 | 0.97518 |
| **3** | **1.18182** | **1** | **0.93223** |
| 4 | 1.18182 | 1.5 | 1.01546 |
| 5 | 1.18182 | 2 | 1.09368 |
| 6 | 1.27273 | 2 | 1.09859 |
| 7 | 1.36364 | 2 | 1.05261 |
| 8 | 1.45455 | 2 | 1.02714 |
| 9 | 1.54545 | 2 | 1.03182 |

Each alternative is simulated 2 times for 10 cycles. The population size is set to 1000 individuals for all alternatives. Column 2 describes the price ratio for X versus Y. Column 3 describes the ratio of ad intensity effort for X versus Y. The values for both these columns are set by the user. The final column describes the ratio of the average purchase motivatios for X versus Y resulting from the simulation of each alternative. The third row (in bold) is the base case where no policy is applied

of a policy Goal. The particular Goal in the Figure has two alternatives depicted as blue nodes and corresponding to the different policy Objectives for reaching the goal and below each Objective lies a yellow node representing the models referred to as Steps that will be used to simulate each alternative. By clicking on any of these nodes the user can see and edit the contents for it in the text area at the bottom of the screen. Using the menu at the top of the screen, she can then update, delete or execute this element or add another element below it. The results of the execution are shown in the text area at the bottom of each screen. For example, the specific scrollable text area in Fig. 3 depicts part of the contents of one of the Objectives indicated with the blue squares in Fig. 3.

We provide a video.[1] that demonstrates the use of the tree-based GUI related to the wine pricing policy. An experimental version of Politika can be found here.[2]

## 3.1 Results and Commentary

We use Table 1 to indicate the types of results and insights that can be reached with Politika using as an example the design of a pricing policy for wine X versus its competitor Y based on the model described in §2. We define our policy goals as a maximum increase in the average purchase motivation for wine X versus Y in the population of interest assuming that we vary the price and ad intensity efforts for X, while leaving unchanged all parameters for Y along with the population parameters in every alternative. As described in §3, the user can create policy alternatives by

[1] http://www.epinoetic.org/Assets/SocSim22.mp4.

[2] http://www.epinoetic.org:4000.

selecting values for policy parameters such as the prices of wines X and/or Y or the ad intensity efforts for X and/or Y.

As Table 1 shows, the best policy alternative among the nine alternatives that were simulated occurs for a price ratio of 1.27273 between X versus its competitor Y and for a doubling of the intensity of the ad effort for X versus Y. At this point we compute the highest purchase motivation for X versus Y (a 1.09859 value for the ratio of the average purchase motivation of X versus Y). This is much better than the base case that estimates the current status between these two wines with no policy applied. For this base case we compute an average purchase motivation for X lower than that of Y (a 0.93223 ratio of average purchase motivations) at a price ratio of 1.18182 , assuming the same intensity of the ad effort between X and Y. Therefore, for the best policy alternative compared to the base case we were able to increase both the purchase motivation for X and the price for X but with a doubling of the intensity of the ad effort for X versus Y. The identification and estimation of such trade-offs between policy variables is one of the advantages of using Politika for simulating and evaluating a series of policy alternatives. This is especially useful for the early, conceptual phase of policy design, during which policymakers need to establish high-level insights for the comparative effectiveness, efficiency, viability and complexity of various alternatives with respect to the policy goals. The outcomes of this phase can then guide strategic decisions on the policy alternatives that should be further pursued and refined in subsequent design stages.

## 4   Conclusions and Related Work

Recently there has been an interest in applying agent-based social simulation in modeling consumer behavior in general [1] and wine purchase motivation in particular [2]. Our efforts complements this research by embedding social simulation in policy design so that the development, analysis and evaluation of alternatives can be facilitated. More information on the inner workings of Politika can be found in [4].

## References

1. Zhang, T., Zhang, D.: Agent-based simulation of consumer purchase decision-making and the decoy effect. J. Bus. Res. **60**, 912–922 (2007). https://doi.org/10.1016/j.jbusres.2007.02.006
2. Huiru, W., et al.: An agent-based modeling and simulation of consumers' purchase behavior for wine consumption. IFAC-Pap On Line **51**(17), 843–848 (2018). https://doi.org/10.1016/j.ifacol.2018.08.089
3. Kyriazis, D. et al.: PolicyCLOUD: analytics as a service facilitating efficient data-driven public policy management. In: Maglogiannis, I. et al. (eds.) Artificial Intelligence Applications and

Innovations (AIAI 2020), IFIP Advances in Information and Communication Technology, vol. 583. Springer, Berlin. https://doi.org/10.1007/978-3-030-49161-1_13

4. Sgouros, N.M.: Politika: implementing a novel meta-simulation methodology for public policy design on the web. Digit Gov Res Pract ACM. https://doi.org/10.1145/3568167

# Exploring Credit Relationship Dynamics in an Interbank Market Benefiting from Blockchain-Based Distributed Trust: Insights from an Agent-Based Model

**Morteza Alaeddini** , **Julie Dugdale** , **Paul Reaidy** , **and Philippe Madiès**

**Abstract** Trust is crucial in economic complex adaptive systems, where agents frequently change the other side of their interactions, which often leads to changes in the system's structure. In such a system, agents who seek as much as possible to build lasting trust relationships for long-term confident interactions with their counterparts decide whom to interact with based on their level of trust in existing partners. A trust crisis refers to the time when the level of trust between agents drops so much that there is no incentive to interact, a situation that ultimately leads to the collapse of the system. This paper presents an agent-based model of the interbank market and evaluates the effects of using a voting-based consensus mechanism embedded in a blockchain-based loan system on maintaining trust between agents and system stability. In this paper, we rely on the fact that blockchain as a distributed system only manages the existing trust and does not create it on its own. Furthermore, this study uses actual blockchain technology in its simulation rather than simply presenting an abstraction.

---

M. Alaeddini (✉) · J. Dugdale
Grenoble Informatics Laboratory (LIG), Université Grenoble Alpes, Grenoble, France
e-mail: Morteza.Alaeddini@univ-grenoble-alpes.fr

J. Dugdale
e-mail: Julie.Dugdale@imag.fr

M. Alaeddini · P. Reaidy · P. Madiès
Grenoble INP, Université Grenoble Alpes, CERAG, 38000 Grenoble, France
e-mail: Paul.Reaidy@univ-grenoble-alpes.fr

P. Madiès
e-mail: Philippe.Madies@grenoble-iae.fr

© The Author(s), under exclusive license to Springer Nature Switzerland AG 2023
F. Squazzoni (ed.), *Advances in Social Simulation*, Springer Proceedings in Complexity,
https://doi.org/10.1007/978-3-031-34920-1_33

# 1 Introduction

Trust seems to be a focus in promoting the ability of agents to collaborate across a complex adaptive system [1, 2]. Models developed in this area seek to avoid the failure of agents' collaborative efforts by entering them into a relationship in order to collaborate [3]. However, trust may be damaged due to unforeseen changes in the environment. In addition to the agents' selfishness due to different ownerships, which sometimes makes them unreliable, one primary source of mistrust in such systems is their lack of global perspective and complete knowledge of the whole environment and their peers with hidden intentions [4].

An interbank market, as a highly stochastic economic environment [5], is a complex adaptive system [6, 7] where banks lend large amounts of money to each other at interbank rates when they need liquidity in a short period [8], thus adapting to this stochastic environment. Establishing more lending relationships in this market helps borrowers with more diverse sources of liquidity [9] and enables them to borrow at lower interest rates from lenders with whom they have a relationship [10]. However, these relationships dynamically change due to the short-term nature of unsecured funding [11]. In order to preserve credit relationships, maintaining a level of trust is essential for all market participants, as its evaporation can lead to instability and liquidity crises [8].

These days, the notion of distributed trust [12] has been reintroduced through the use of blockchain [13]. As a cryptographically secured, distributed ledger, this technology is widely believed to spread trust in digital environments [14]. In this study, using the aggregate balance sheet of French banks, we model an interbank market as a multi-agent system and examine whether blockchain is able to compensate for the loss of trust among peers during economic declines. Concretely, the contributions of this paper are twofold: (i) adding to the literature on trust in multi-agent systems and (ii) using blockchain as part of the simulation platform. The rest of the paper is organized as follows: Sect. 2 gives background information on related work previously performed in this area. Section 3 describes the components of the model and the behavior of various agents in different circumstances. The results of simulating this model based on a number of scenarios are presented in Sect. 4. Finally, Sect. 5 concludes the paper and gives avenues for future research.

# 2 Related Work

## 2.1 Distributed Trust

The notion of distributed trust is not new and dates back to the late 1990s [12]. Among the methods proposed for building trust in multi-agent systems, one can find those that benefit from this notion. Jordi and Sierra [15] use a reputation mechanism in which each agent records its direct trust in other agents resulting from interacting

with them in a local database and shares these data with other agents so that they use them in their indirect trust estimation. Jurca and Faltings [16] propose a set of broker agents responsible for gathering reports from other agents on their interactions with each other. The broker agents also provide reputation information to agents who need it. Tweedale and Cutler [17] attribute trust to the collective decision of a hierarchical team of which the agent is a member. Huynh et al. [18] integrate all of these methods into a framework called FIRE. However, in the past, there were many obstacles to the objectification and implementation of distributed trust in practice because it is unreasonable to expect such information to be shared by all members of the system [18].

By using blockchain, which refers to a cryptographically secured distributed ledger with a decentralized consensus mechanism, it is easier to implement such ideas. Calvaresi et al. [19] provide a JADE-based architecture and implement a system that computes agents' reputations using smart contracts and enables tracking of how their reputation changes. Khalid et al. [20] propose maintaining trust in an agent-based distributed energy trading system by publishing information on inter-agent agreements in the blockchain. Alaeddini et al. [21] consider blockchain in designing a multi-agent interbank trading system, where trust is regarded as a significant concern. It is worth noting that none of these studies addresses an individual agent's threshold for the trust it needs to have in another agent to interact, and in fact, they all have given the same recommendation to all agents, regardless of their specific characteristics. Also, none of the models uses a real blockchain as part of their simulation system.

## 2.2 The Selected Trust Mechanism

Unlike the mentioned methods of trust in multi-agent systems, we propose a new method to develop a trust model based on the consensus reached by agents and using some variables found by Bülbül [22]. The following features are the main distinctions of this method from others:

- Both the expected level of trust of the agent responding to the interaction and the level of trust met by the agent requesting the interaction are considered;
- It uses a blockchain-based consensus algorithm to establish distributed trust; and
- Unlike some other methods (e.g., Khalid et al. [20]), it does not publish any confidential information of agents on the blockchain.

The model uses a reputation system as an additional trust layer based on counterparts' relationships [23] and applies six values from $-1$ (distrust) to 4 (complete trust) for both direct and indirect trust. The value of direct trust is the result of assessing the lender's trust in a borrower for a loan transaction, while an indirect value is based on reputation information. An agent uses values of its direct trust in other agents in order to arrive at a consensus on their reputation and recommend them to other agents.

To calculate the level of trust desired by the lender agents, the model uses three determinants, including current interaction with the central bank ($X_{1i,t}$), equity ($X_{2i,t}$), and size ($X_{3i,t}$) of the lending bank as follows at every time step $t = 1, 2, \ldots, \text{T}$:

$$X_{1i,t} = \begin{cases} 1, & \text{bank } i \text{ has a debt to the central bank in time } t \\ 0, & \text{otherwise} \end{cases} \tag{1}$$

$$X_{2i,t} = \begin{cases} 2, & \text{bank } i'\text{s capital at time } t > \text{average capital of banks of similar size} \\ 1, & \text{bank } i'\text{s capital at time } t = \text{average capital of banks of similar size} \\ 0, & \text{otherwise} \end{cases} \tag{2}$$

$$X_{3i,t} = \begin{cases} 1, & \text{bank } i \text{ is small in size} \\ 0, & \text{bank } i \text{ is medium in size} \\ -1, & \text{otherwise} \end{cases} \tag{3}$$

$$\tau_{i,t} = \max\left[-1, \left(X_{1i,t} + X_{2i,t} + X_{3i,t}\right)\right] \tag{4}$$

where $\tau_{i,t}$ denotes bank $i$'s observed ordinal variable as the trust threshold at time $t$.

Let $\theta_{i,j,t}$ denote the trust level between banks $i$ and $j$ at time $t$. A lending relationship between lender $i$ and borrower $j$ is allowed at time $t$ if $\theta_{i,j,t} + \tau_{i,t} > 4$. The level of direct trust between banks $i$ and $j$ at time $t$ is obtained from the Eq. 5, where $H_{i,j,t+}$ indicates the history of good records of bank $j$ in repaying the loans it has received from bank $i$ until time $t$, and $H_{i,j,t-}$ denotes the history of bad records in the same period. $H_{i,j,t}$ indicates the number of lending relationships between banks $i$ and $j$ from the beginning to time $t$.

$$\theta_{i,j,t} \approx \max\left[-1, \left(\frac{H_{i,j,t+} - H_{i,j,t-}}{H_{i,j,t}} \times 4\right)\right] \tag{5}$$

To calculate the level of indirect trust between two peers who did not have a lending relationship with each other before time $t$, each agent asks endorsers to examine the level of trust of the new counterpart. Endorsers, which are nodes located in the credit paths leading to the endorsee node, calculate its trust level by tracking the points assigned to that node and their credit paths that lead to it. The final score, subject to consensus, would be obtained based on the recommendations of other related nodes before the time of the loan transaction. The model defines the indirect trust of bank $i$ to bank $j$ as the direct/indirect trust of the counterparts $k$ of bank $i$ to bank $j$, weighted by the trust of bank $i$ towards these neighbor agents $k$. Let $w_{i,j,t}$, defined as follows, denote the elements of the stochastic matrix for normalizing the values of $\theta_{i,j,t}$ (= 0 if there is no link between agents $i$ and $j$).

$$w_{i,j,t} = \frac{1}{4} \times \frac{\theta_{i,j,t}}{\text{n}\left(N_{i,t}\right)} \tag{6}$$

where $N_{i,t}$ is the set of neighbors of agent $i$ at time $t$, and $n(N_{i,t})$ denotes the number of elements in this set. The indirect trust score of bank $i$ to bank $j$ is calculated as follows:

$$\theta_{i,j,t} \approx \sum_{k \in N_{i,t}} w_{i,k,t} \theta_{k,j,t} \qquad (7)$$

This means that in order to calculate the level of trust of a counterpart if there is a direct relationship, agents use (5); otherwise, they need the consensus of other agents based on (7) (maybe in a recursive mode).

## 3   The Model

### 3.1   The Agent-Based Simulation System

Our model developed in Repast Simphony builds on a number of recent studies [5, 24–26] and is populated by two types of agents: (i) N banks that interact with and lend to each other, and (ii) one central bank that regulates the market and helps banks avoid failure when necessary. Furthermore, a lending contract has been developed to support interactions among the agents.

Banks are heterogeneous, imperfect, autonomous, and, to some degree, adaptive agents. They follow base-level rules to make interbank placements and must meet all regulatory requirements in their transactions and changes in their balance sheets. The initiation stage in our model creates random counterparts for banks, assigns their initial assets and liabilities according to banks' sizes, and determines each bank's balance sheet. The natural and financial sides of the market are linked by multiple, non-linear feedbacks and evolve in a finite time horizon. In each time step (one day), the items on the banks' balance sheets change stochastically by following Gaussian random walks with related moving drifts $\mu_{o_{i,t-1}}$ and noises $\sigma_{o_i}$ (see Sect. 3.2 for details).

The general logic of the simulation is that banks manage their liquidity (cash) by exchanging funds in the market. It is assumed that, at first, there are no loans to be repaid by banks (none of the banks owes to other banks). After the change in the banks' balance sheets in the first time step, the interaction of banks to borrow funds overnight in order to compensate for their lack of liquidity forms the interbank lending network in our model. The payments settlement is managed by a central clearing counterparty (i.e., the central bank), and all interbank loans are simulated to be paid in the blockchain (see Sect. 3.2). Figure 1 shows the sequence of actions performed at each time step.

As shown in Fig. 1, at the beginning of each period, the amounts of clients' deposits, loans, and interbank payments resulting from the total transactions of clients with clients of other banks are updated stochastically. The central bank makes

**Fig. 1** The simulation process in BPMN (Y: yes; N: no; C: compensated; UC: uncompensated)

a clearing matrix for the payments, and banks use their reserve balance to settle their clearing vector. Then banks repay their matured interbank debts by their cash (reserve) balance and are evaluated by the lender. Banks that do not have enough reserves to repay their debts, if they have credit receivable on the same day due to the repayment of other banks' debts, wait until the successive settlement cycles on the present period; otherwise, they repay their debts by borrowing first from their counterparts and then from other banks (see Fig. 2). Banks then calculate their liquidity excess or deficit and provision their reserve. Banks that have excess liquidity pay part of the surplus to buy securities (investment) and then lend to other banks, according to Fig. 2. Finally, if banks owe money to the central bank, they repay it.

As shown in Fig. 2, in order to manage their liquidity, banks borrow or lend in the market. For this purpose, based on their history, borrowing banks send their loan requests first to their lending counterparts. Lending partners respond to requests based on their excess and borrowing banks' history. If banks cannot borrow from their existing counterparts, they will apply for a loan from other banks with a lending position in the network. If lending banks meet all or part of the liquidity needs of the applicant banks, they adopt two different strategies against the two off-chain (traditional) and on-chain (blockchain-based) modes (see Sect. 4.1). In either case, both borrowing and lending banks add each other to their counterparty list if the loan is agreed upon.

Banks that have not been able to make up for their need in the market will be refinanced by the central bank if they have enough securities; otherwise, they will have to fire sell—selling assets at heavily discounted prices. Then, banks try to repay
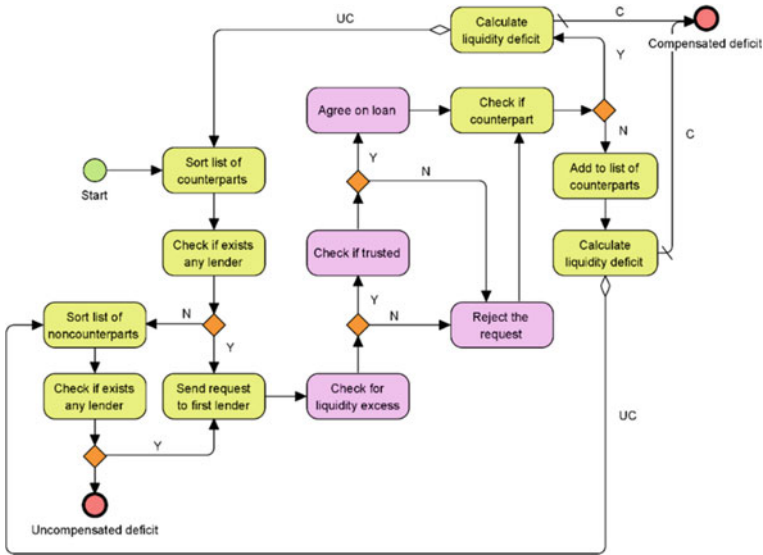
**Fig. 2** The loan process in BPMN (Y: yes; N: no; C: compensated; UC: uncompensated)

their overdue loans, if relevant. At the end of each period, a bank goes bankrupt if it fails to make up for its liquidity deficit or its equity is zero or less and does not compensate for these problems by raising its equity. The failed bank is removed from the model. The bank's failure also leads to losses resulting from its zero debt to the banks from whom it has borrowed. This is the unique source of systemic risk and instability in our model. It is worth noting that the flow diagrams in Figs. 1 and 2 represent real bank behaviors [24, 25].

## 3.2   The Blockchain-Based Loan System

Adapted from Cucari et al. [27], we develop a simple loan system on the consortium blockchain Corda that records loan transactions of agents. To develop this system, we use the logic of a simple CorDapp already implemented by the Corda team and make changes based on our specific needs. One of the items that the loan system records and maintains is loan state, which is an immutable object representing facts (loan data) known only by counterparts. The system also benefits from smart contracts between banks by turning the contract terms into code that executes automatically when they are met. The contract code is replicated on the nodes in the network. All these nodes have to reach a consensus that the terms of the agreement have been met before they execute the contract. Figure 3 shows the sequence of consensus in the system.
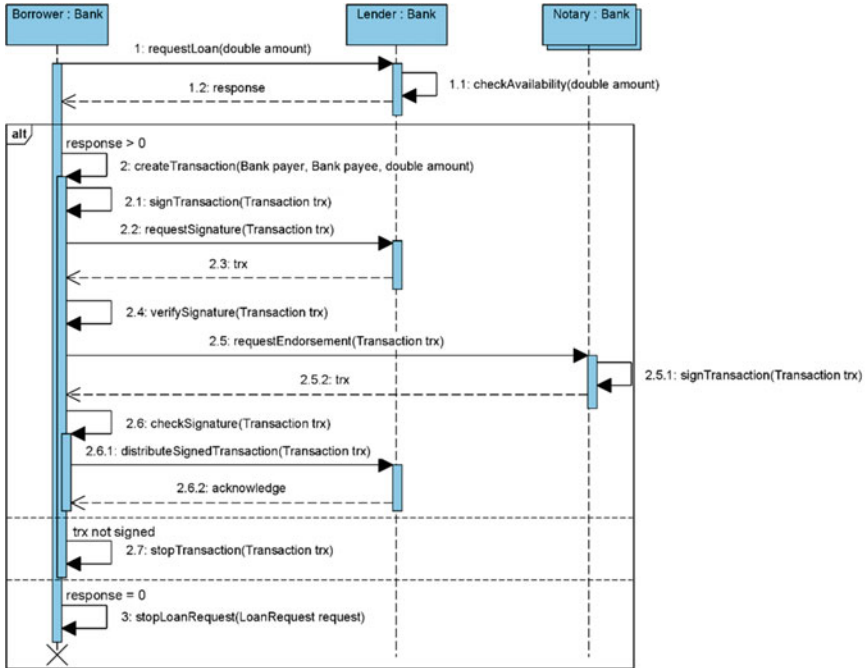
**Fig. 3** The UML sequence diagram of consensus in the blockchain-based loan system

Loan transactions must achieve both validity and uniqueness consensus to be committed to the ledger. The first determines if a transaction is accepted by the smart contracts it references, while the latter prevents double-spends, i.e., the risk that the money is paid twice or more. A transaction must have all the necessary signatures to reach the validity consensus, meaning that the qualification of a borrower who has no previous relationship with a particular lender must be endorsed by a notary consisting of the banks that have already lent to that borrower. Uniqueness consensus is when the notary checks that the lender has not used the same input for multiple transactions.

Communication between banks is point-to-point using a flow, which automates the process of agreeing on ledger updates between the banks. Our agent-based simulator communicates with the loan system through an API that we developed. The initiation stage in our model deploys one node in the blockchain for each agent. The deployed blockchain nodes containing the API that records loan transactions on the blockchain are then run at this stage. Therefore, the environment we implement to simulate agents' behavior is as similar as possible to the real environment that banks may use in a real market by employing a real blockchain to record their loan transactions.

**Table 1** Parameters for different economic cycles

| Parameter | Growth | Decline | Recession |
|---|---|---|---|
| Noise of credits and lending | $\mathcal{U}(0, 0.005)$ | $\mathcal{U}(0.05, 0.1)$ | $\mathcal{U}(0.1, 0.25)$ |
| Noise of deposits and payments | $\mathcal{U}(0, 0.003)$ | $\mathcal{U}(0.03, 0.06)$ | $\mathcal{U}(0.06, 0.15)$ |

## 4 Experimental Evaluation

### 4.1 Scenarios

We study the interbank market dynamics with and without using blockchain. We first test three economic cycle scenarios in the absence of blockchain (off-chain mode) using the parameters from a uniform distribution shown in Table 1. In the next step, we intervene with consensus in the blockchain on the level of trust between banks and test the three scenarios again (on-chain mode). To investigate the number of simulations required to smooth out irregularities, we apply the convergence of subsequent mean values at the aggregation level by forming a moving mean value. As soon as the deviation of the calculated mean value from the convergence mean value is less than 0.05, we consider it to be robust. Although 40 simulations on average are enough to reach a robust mean, we only run each simulation ten times because of time constraints. Finally, we compare the average of results of these six experiments.

Each of these setups assumes that banks face an abundance or lack of liquidity with specific dynamics. The values of the other parameters used in our study are the same for all scenarios and are according to the coefficients and minimums set in Basel III and enforced by the ECB.
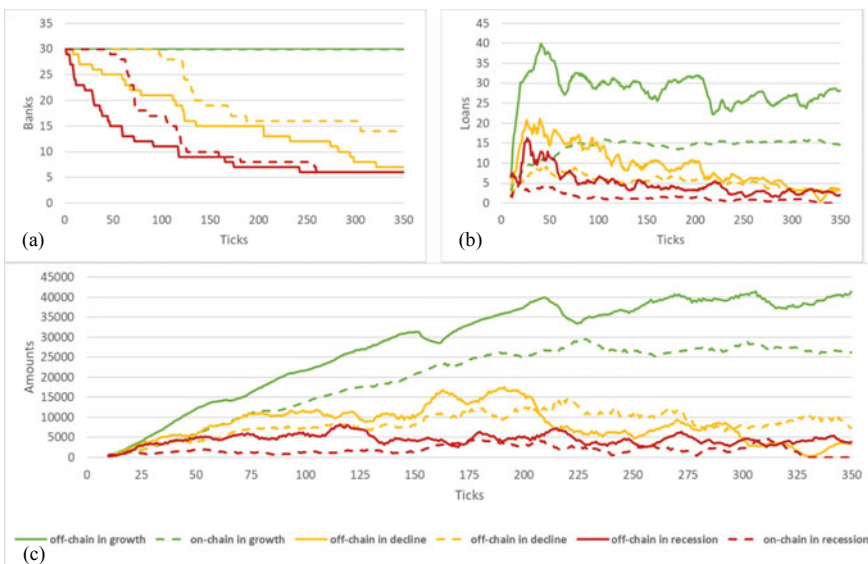
### 4.2 Experiments

Due to the limited hardware resources for simultaneous allocation to the blockchain network nodes, we perform this simulation with 30 random banking agents of different sizes whose balance sheets are adjusted based on the characteristics of banks operating in the French interbank market. Each experiment is based on an independent draw of the initial network for interbank loans as well as the balance sheet amounts of different banks. All initial networks are formed based on an initial value of 3 as the number of each bank's counterparts. However, as the simulation progresses, it is possible for banks to make new relationships over time. In off-chain mode scenarios, the acceptance of the request by a new lender is $1 - \sigma_{l_i}$ likely but at a higher premium than what the other partners of that agent pay. In on-chain mode

scenarios, conditions in Sect. 2.2 must be met for the request to be accepted by a
new lender, and the interest rate will be subject to the same procedure as the agent's
other counterparts.

Comparing the number of banks in both off-chain and on-chain modes during
350 days of activity in times of economic growth in Fig. 4a shows the stability of
banks throughout the period. This means that when uncertainty is low, banks meet
each other's liquidity needs well, and no default or failure is observed. From Fig. 4b,
as the economy grows, the total number of loans in 40% of days in the off-chain
mode is in a position above the 'number of banks' line ($n = 30$). It seems that in
this mode, the general tendency of the borrowing banks is to borrow from more
partners, and most of the lending banks tend to diversify their loan portfolio and lend
to more banks. Another possible reason for this could be the dual activity of some
banks, which act as both the lender and borrower in one day. By the intervention
of blockchain, this trend goes surprisingly below the line (100%), meaning that a
higher level of trust is interpreted as having narrower yet a deeper relationships with
peers.

A look at the starting points and progression trends of the failures in Fig. 4a
provides a similar comparison for the economic decline situation. The simulation
results of the employed consensus algorithm to build trust among market members
show that although the use of blockchain to realize this algorithm cannot ultimately
prevent cascades of banks' failure, they can delay the collapse at least for a consid-
erable time (89 days). During a recession, this opportunity is reduced to 45 business
days. This means that the impact of using blockchain in restoring trust in the market



**Fig. 4** Evolution of banks (**a**), the number of interbank loans (**b**), and total amount of interbank
loans (**c**) in times of different economic situations: off-chain versus on-chain modes

during an economic decline is almost twice as great as during a recession. According to the almost uniform distance between the two on-chain diagrams in these two states in Fig. 4a, it seems that this is more affected by the severity of uncertainty caused by the economic situation than due to the blockchain features.

Similar to economic growth scenarios, there are significant gaps between off-chain and on-chain modes in both the number and the total amount of loans in times of economic decline and recession (see Fig. 4b and c). It means that if the blockchain is used in an unstable economic situation, both parties would prefer to reduce the number of agreements and increase the amount instead (i.e., trustworthy relationships). The gap between the number of loans in these two modes remains almost constant.

## 5 Conclusion and Future Work

This paper presents a multi-agent simulation platform for the interbank market and integrates the notion of trust using a blockchain-based consensus algorithm to explore dynamics of lending relationships and the effects of uncertainty caused by different economic conditions. In order to compensate for the lack of liquidity of their peers, banks lend each other overnight. Unlike Khalid et al. [20], we do not publish information on inter-agent agreements in the blockchain. Instead, we ask endorsing nodes in the network to determine how reliable is the borrower based on their previous track records, and finally to validate the transaction through a voting mechanism.

Based on the simulation results, the banking network remains stable during periods of economic growth without any additional need for a mechanism to strengthen trust. The issue of which of these two strategies in times of economic growth leads to lower cost and more operational advantage for the system and members can be the subject of new research. However, by increasing the uncertainty caused by changes in economic conditions, the establishment of a blockchain-based consensus mechanism in the market can help maintain trust between banks and, consequently, system stability (i.e., continuation of the presence of agents in the system). Although such a mechanism is not able to fully protect the market from contagious failures in the long run, it undermines the destructive effects of uncertainty for a significant period. An important point for the regulator and market participants is that since blockchain is an important factor in ensuring market resilience, the resiliency of the blockchain infrastructure should also be taken into account in times of stress so that it can meet expectations.

Good and bad history kept by the agents in our model can be interpreted as belief and disbelief [28]. Also, because uncertainty is considered as a parameter affecting the agents' behavior, we are interested in combining our method with the method of Cheng et al. [28], which basically uses these items in calculating trust. As a limitation of our model, banks' decisions about lending, like other events outside their control, have a stochastic basis. Another future agenda is adding learning capabilities to the model so that agents make decisions based on their current and future goals, use

what they learn from the past, and consider other agents' behavior. Maintaining trust between agents can be one of the goals to which they apply what they learn in using blockchain to conduct more trustful transactions in the future. Furthermore, the results do not model the case of a black swan event that could be the cause of a systemic collapse. A scenario in which economic growth is abruptly followed by recession can be of interest to scholars and practitioners to analyze the market in off-chain and on-chain modes.

# References

1. Ramchurn, S.D., Huynh, D., Jennings, N.R.: Trust in multi-agent systems. Knowl. Eng. Rev. **19**(1), 1–25 (2004)
2. Kim, W.-S.: Effects of a trust mechanism on complex adaptive supply networks: an agent-based social simulation study. J. Artif. Soc. Soc. Simul. **12**(3), 4 (2009)
3. Wooldridge, M., Jennings, N.R.: The cooperative problem-solving process. J. Log. Comput. **9**(4), 563–592 (1999)
4. Pinyol, I., Sabater-Mir, J.: Computational trust and reputation models for open multi-agent systems: a review. Artif. Intell. Rev. **40**(1), 1–25 (2013)
5. Steinbacher, M., Jagrič, T.: Interbank rules during economic declines: can banks safeguard capital base? J. Econ. Interac. Coord. **15**(2), 471–499 (2020)
6. Chiriță, N., Nica, I., Popescu, M.: The impact of bitcoin in the financial market. A cybernetics approach. In: Education, Research and Business Technologies. Springer, pp. 197–209 (2022)
7. Glass, R.J., Ames, A.L., Brown, T.J., Maffitt, S.L., Beyeler, W.E., Finley, P.D., Linebarger, J.M.: Complex adaptive systems of systems (CASoS) engineering: mapping aspirations to problem solutions. Sandia National Laboratories, Albuquerque, New Mexico, USA (2011)
8. Alaeddini, M., Madiès, P., Reaidy, P., Dugdale, J.: Interbank money market concerns and actors' strategies—a systematic review of 21st century literature. J. Econ. Surv. **37**(2), 573–654 (2023)
9. Craig, B.R., Fecht, F., Tumer-Alkan, G.: The role of interbank relationships and liquidity needs. J. Bank. Financ. **53**, 99–111 (2015)
10. Cocco, J.F., Gomes, F.J., Martins, N.C.: Lending relationships in the interbank market. J. Financ. Intermediation **18**(1), 24–48 (2009)
11. Anand, K., Gai, P., Marsili, M.: Rollover risk, network structure and systemic financial crises. J. Econ. Dyn. Control **36**(8), 1088 (2012)
12. Abdul-Rahman, A., Hailes, S.: A distributed trust model. In: Proceedings of the 1997 Workshop on New Security Paradigms, pp. 48–60 (1998)
13. Calvaresi, D., Dubovitskaya, A., Calbimonte, J.P., Taveter, K., Schumacher, M.: Multi-agent systems and blockchain: results from a systematic literature review. In: International Conference on Practical Applications of Agents and Multi-agent Systems. Springer, pp. 110–126 (2018)
14. Shin, D.D.: Blockchain: the emerging technology of digital trust. Telematics Inform. **45**, 101278 (2019)
15. Jordi, S., Sierra, C.: Regret: a reputation model for gregarious societies. In: Proceedings of the Fourth Workshop on Deception Fraud and Trust in Agent Societies, Montreal, Canada, pp. 61–70 (2001)

16. Jurca, R., Faltings, B.: Towards incentive-compatible reputation management. In: Trust, Reputation, and Security: Theories and Practice: AAMAS 2002 International Workshop, Bologna, Italy, 15 July 2002. Selected and Invited Papers. Springer, p. 138 (2003)
17. Tweedale, J., Cutler, P.: Trust in multi-agent systems. In: International Conference on Knowledge-Based and Intelligent Information and Engineering Systems. Springer, pp. 479–485 (2006)
18. Huynh, T.D., Jennings, N.R., Shadbolt, N.R.: An integrated trust and reputation model for open multi-agent systems. Auton. Agent Multi-agent Syst. **13**(2), 119–154 (2006)
19. Calvaresi, D., Mattioli, V., Dubovitskaya, A., Dragoni, A.F., Schumacher, M.: Reputation management in multi-agent systems using permissioned blockchain technology. In: 2018 IEEE/WIC/ACM International Conference on Web Intelligence (WI). IEEE, pp. 719–725 (2018)
20. Khalid, R., Samuel, O., Javaid, N., Aldegheishem, A., Shafiq, M., Alrajeh, N.: A secure trust method for multi-agent system in smart grids using blockchain. IEEE Access **9**, 59848–59859 (2021)
21. Alaeddini, M., Dugdale, J., Reaidy, P., Madiès, P., Gürcan, Ö.: An agent-oriented, blockchain-based design of the interbank money market trading system. In: Agents and Multi-Agent Systems: Technologies and Applications 2021. Springer, pp. 3–16 (2021)
22. Bülbül, D.: Determinants of trust in banking networks. J. Econ. Behav. Organ. **85**, 236–248 (2013)
23. Pazaitis, A., De Filippi, P., Kostakis, V.: Blockchain and value systems in the sharing economy: the illustrative case of backfeed. Technol. Forecast. Soc. Change **125**, 105–115 (2017)
24. Popoyan, L., Napoletano, M., Roventini, A.: Winter is possibly not coming: mitigating financial instability in an agent-based model with interbank market. J. Econ. Dyn. Control **117** (2020)
25. Sato, A.H., Tasca, P., Isogai, T.: Dynamic interaction between asset prices and bank behavior: a systemic risk perspective. Comput. Econ. **54**(4), 1505–1537 (2019)
26. Teply, P., Klinger, T.: Agent-based modeling of systemic risk in the European banking sector. J. Econ. Interac. Coord. **14**(4), 811–833 (2019)
27. Cucari, N., Lagasio, V., Lia, G., Torriero, C.: The impact of blockchain in banking processes: the Interbank Spunta case study. Technol. Anal. Strateg. Manage. **34**(2), 138–150 (2022)
28. Cheng, M., Yin, C., Zhang, J., Nazarian, S., Deshmukh, J., Bogdan, P.: A general trust framework for multi-agent systems. In: Proceedings of the 20th International Conference on Autonomous Agents and MultiAgent Systems, pp. 332–340 (2021)

# Simulating Bounded Rationality in Decision-Making: An Agent-Based Choice Modelling of Vehicle Purchasing

**Khoa Nguyen, Valentino Piana, and René Schumann**

**Abstract** This paper investigates the possibility of simulating bounded rationality effects in an agent's decision-making scheme by limiting its capability of perceiving information and utilising a decision-making framework of Triandis' Theory of Interpersonal Behaviour. Based on previous work on an agent-based platform, BedDeM, we propose how to capture the effects of sequential, emotional, habitual and multi-criteria decision-making. The Perception component in the agent is further extended to take into account confirmation bias and the bandwagon effect. We demonstrate the functionality of this model in the context of purchasing vehicles in Switzerland's households.

**Keywords** Agent-based simulation · Bounded rationality · Choice modelling · Behavioural theory

## 1 Introduction

The number of agent-based models (ABM) used to represent human decision making is increasing. Agent designs with notion of perfectly rational maximise expected utility but crucially ignore the resource costs incurred. Researches in bounded rationality (BR) offer an alternative to how to model behaviours in an uncertain environment with limited available cognitive resources. However, the ABMs utilised in these researches often focus only on simulating one particular type of BR (see surveys such as [6, 15] and Sect. 2). This study looks at an universal approach of devel-

K. Nguyen (✉) · V. Piana · R. Schumann
HES-SO Valais Wallis, SILab, Rue de Technpôle 3, Sierre 3960, Switzerland
e-mail: khoa.nguyen@helmholtz-hzi.de

V. Piana
e-mail: valentino.piana@hevs.ch

R. Schumann
e-mail: rene.schumann@hevs.ch

oping an agent-based platform that can investigate the impact of multiple BRs on decision-making.

Discussing the term *bounded rationality* equals walking on a tightrope due to different interpretations across and even within disciplines. In this study, we follow the definition provided by Carley et al. [5] regarding two types of bounds in agents—limits to capabilities and limits to knowledge. Capabilities are related to the agent's physical, cognitive and computational architecture. Knowledge is the ability to learn and construct intellectual history. This paper attempts to take advantage of active perception to limit the agent's capability to observe relevant information. Through this data filtering capacity, BR is an extension of the model of the perfectly informed, optimised individuals to account for restricted knowledge and resources, i.e. a form of *bounded optimality* [22] [p. 1050]. Coupling this definition with the notions of bounded rationality coined by Simon [23] and the heuristics and biases advanced by other researchers, several phenomena can be targeted in this study:

- *Sequential decision-making* refers to algorithms that consider the dynamics of the world, thus delaying parts of the problem until they must be solved [8, 337].
- *Emotional decisions* happen when the people's emotional state influences the depth of information processing related to decision-making [24].
- *Habit formation* is the process by which a behaviour becomes automatic when it is repeated with a routine [24].
- *Multiple criteria* other than cost can be considered, depending on the decision-making context [22][p. 622–628].
- *Confirmation bias* is the tendency of people to select the information that supports their views, ignore contrary information, or when interpret ambiguous evidence as supporting their existing beliefs or values [18].
- *Bandwagon effect* is a psychological phenomenon in which an idea or belief is being followed because everyone seems to be doing so [14].

We acknowledge that this list is limited and only covers the general ideas of each BR. However, it represents topics that are often mentioned in ABM research (see surveys such as [6, 15]) and provides a starting point for what can be considered in our study.

Previously, we have developed an agent-based model, and integrated tooling—BedDeM—based on Triandis' Theory of Interpersonal Behaviour (TIB) [16, 17]. The decision-making modules in this model can be used to implement different mechanisms representing items from the list above. In particular, it currently factors in the effects of *sequential*, *emotional*, *habitual* and *multiple criteria* decision-making (see Sect. 3.4). We modify the *Perception* component to take into account *confirmation bias* and *bandwagon effect*.

The business of purchasing new vehicles is an essential field for Switzerland's energy strategy, especially when it provides an understanding of the need of individual consumers and requirements for future infrastructure [4]. It is also an area where BRs are particularly pervasive, as decisions are made at the level of deeply heterogeneous individuals and households. Due to the significant number of individual decision-makers involved and alternatives offered in vehicle purchasing, ABMs are often

utilised for the assessment of BR's effects in the lab as well as in the field (e.g. [10, 13]. Therefore, it is chosen as a suitable context to implement and test the functionality of the new bounded *Perception* component.

The paper is organised as follows: After considering some of the related ABM architectures in Sect. 2, we present the structure of our agent-based model and explain how the mentioned BRs are specified in Sect. 3. Next, a case study is provided to evaluate the result of applying this bounded Perception in Sect. 4. Finally, we conclude and suggest further development in Sect. 5.

## 2 Related Works

This section provides the state-of-the-art in terms of ABM that addressing the BRs mentioned, i.e. sequential, emotional, habitual, multi-criteria decision-making, confirmation bias, and the bandwagon effect. Our agent decision-making architecture, which also covers several different types of BRs, will be discussed in Sect. 3.

In terms of sequential decision-making, researchers in ABM often take the approach of multiple steps/stages in decision-making before the final output. The most famous architecture of this category is Belief-Desires-Intention (BDI) model [9]. It is centred around three mental attitudes, namely beliefs (the informational state of the agent), desires (the objectives or situations that the agent would like to accomplish or bring about) and, especially, intentions (the deliberative state of the agent - what the agent has chosen to do). Other extensions of BDI, cognitive and normative architectures that have a perception-deliberation-action cycle also belong to this category. A good summary of them can be found in [2].

There is a body of work focusing on emotions in BDI agent reasoning (see [2]). However, only a few agent architectures considered emotions explicitly in literature. These include PECS [27], Emotional BDI (eBDI) [19] and BRIDGE [7]. The first of these is an extension of the BDI architecture that incorporates emotions as a decision criterion into the agent's decision-making process. PECS aims to enable integrative modelling of physical, emotional, cognitive and social influences within a component-oriented agent architecture. BRIDGE represents emotions using the *Ego* component to specify different emotional responses to various stimuli. According to [2] and the best of our knowledge, these architectures are used as reference models, so few specifics can be found about their actual implementations in practice.

To represent the habitual patterns in human behaviour, hybrid approaches that allow for heuristics, as well as deliberation and reactive production rules, are often utilised in ABM. Two examples of this category are Consumat [12] and BRIDGE [7]. Consumat allowed for modelling habitual behaviour by introducing five heuristics based on uncertainty and cognitive effort that can be utilised instead of complete deliberation. BRIDGE, similar to Consumat, introduces the idea of the basic needs of the agent, which can overrule any deliberate decision-making process via a response component to ensure that agents can react when needed.

Multi-criteria decision-making is usually addressed by applying *multi-attribute utility theory*, which is used to represent the preferences of an agent over bundles of goods either under conditions of certainty about the results of any potential choice or under conditions of uncertainty [22][p. 622–628]. To consider the attribute that is not mutually utility independent, Thiriot et al. also propose a multi-objective multi-agent system (MOMAS) to explicitly consider the possible trade-offs between conflicting objective functions [26]. The criteria are often context-dependent, i.e. the modeller has to define them based on statistics or previous empirical studies.

Confirmation bias considers how various sources of information are filtered due to personal cognitive biases. For example, eBDI filters information from all perceptions and other sensor stimuli using semantic association rules derived from its internal beliefs. BRIDGE architecture has an *Ego* component that contains different filters and ordering preferences. They are utilised to interpret the input stream of information to form the beliefs in the agent. Confirmation bias is also considered under the opinion dynamics modelling frameworks. Sobkowicz introduced a quasi-Bayesian belief updating framework, where the incoming information is filtered by the cognitive biases or predispositions of the agent (e.g. memory priming/availability, simplicity/attention and emotional filters) [25]. Rollwage et al. suggest implementing confirmation bias via meta-cognition (accuracy of belief formation) of agents, allowing them to down weight contradictory information when correct but still able to seek new information when they realise they are wrong [21].

The bandwagon effect can be associated with the ability to consider social learning in agent design, which is often found in normative models. Several architectures can be listed in this category, including BRIDGE, EMIL-A [1] and Consumat. BRIDGE accounts for some social concepts, including a social interaction consideration, the social concept of culture, and a notion of self-awareness (and resulting differentiation of one-self and other agents). In EMIL-A, social norms instead play a central role. It models the process of agents learning about norms in a society, the internalisation of norms and the use of these norms in the agents' decision making. On the social level, Consumat has some idea of sociality in terms of agents being able to reason about the success of their actions in relation to the success resulting from the actions of their peers. If the agent does not perform as well, it simply imitates (i.e. copies) the action(s) of others.

Although some account for multiple aspects of behaviour, the agent architectures and implementations surveyed above do not comprehensively cover all BR effects mentioned in Sect. 1. Therefore, in this study, we create an agent model capable of considering these effects in its decision-making scheme.
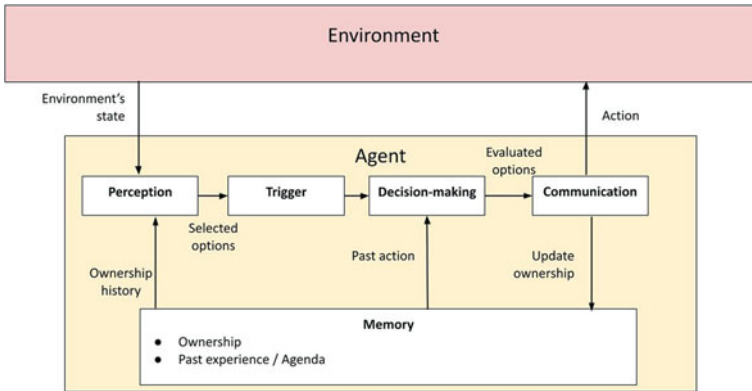
**Fig. 1** Overview of agent's components

# 3 Simulating Bounded Rationalities in Agent's Decision-Making

Several effects of BRs can already be covered using our previous work on an agent framework based on TIB [17], including sequential, emotional, habitual and multi-criteria decision-making. We further extend the *Perception* component to cover the confirmation bias and the bandwagon effect. In the first subsection, we provide an overview of the agent's decision-making cycle. The following subsections describe the two main components related to this study: *Perception* and *Decision-making*. We then summarise how each BP type has been captured in our agent architecture.

## 3.1 Agent's Decision-Making Cycle

The main components of our agent's decision-making cycle are illustrated in Fig. 1. First, it uses the *Perception* component to observe information about the available options. Using the agent's reference, it then filters, sorts, and creates a shortlist of options. If the agent's internal state or these options satisfy specific criteria, the *Decision-making* component is triggered. It follows the procedure of the TIB framework to evaluate the list of options in terms of a utility value (detailed below). Finally, an option is selected based on the provided utility, either by choosing the best (deterministic agent) or using a probability (probabilistic agent). The *Communication* component then outputs this action to the environment and updates the *Memory* component of the agent.
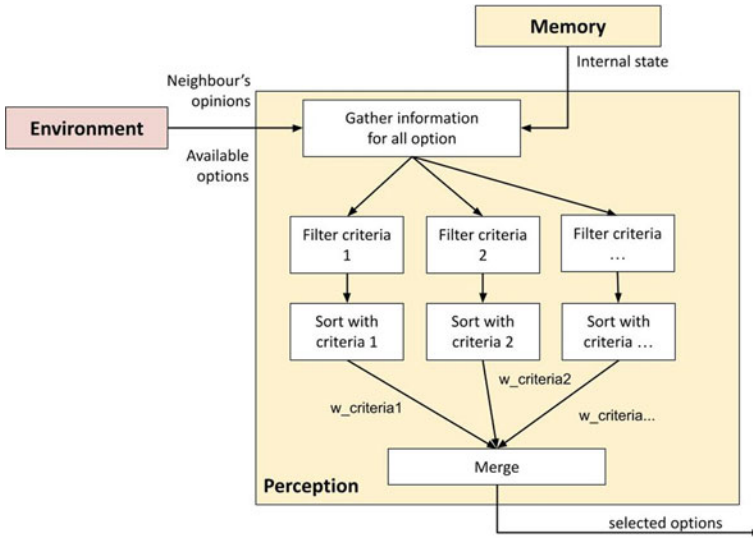
**Fig. 2** Perception component

## 3.2 Perception Component

The *Perception* component (see Fig. 2) first gathers information about the available options from the environment, including its neighbour's opinion. It then divides them into several lists, each satisfying certain criteria. These lists are then sorted, multiplied with certain weights and merged to form a list of selected options for decision-making. The criteria and their weights are based on the agent's personal preferences about the option's properties, which can be calibrated using empirical data.

The mechanism can be explained clearer in the context of car purchasing: a consumer often starts by filtering out models that have a certain type of engine, price, energy labels and neighbour reviews. As human mental accounting mechanisms are limited [11], s/he has to sort the options to get the best one from each category and combine them to make a final list of available models for the ultimate decision-making step.

Using this structure, the confirmation bias can be represented through the filtering process with only relevant options being considered. The bandwagon effect is highlighted with the inclusion of neighbour's opinion as one of the criteria. Using an associated weight, the agent can decide on the influence of this effect on its final list of selected options.

## 3.3  Decision-Making Component

A full decision-making component with the TIB framework is illustrated in Fig. 3.
For all determinants ($d$), each option ($opt$) is given a utility value which comes
from comparing its property with other's ($U_d(opt)$). In the first level, this value can
be in the form of a real value numerical system (for determinants such as price or
time) or ranking function (for determinants such as emotion). Either of which can
be calculated from empirical data (e.g. census, survey) or calibrated with expert
knowledge and stakeholders' assessment. The results for these determinants are then
normalised and multiplied with an associated weight (called $w_d$); the sum of which
becomes the reference value for the option at the next level. This process is captured
in the following equation:

$$EU_d(opt) = \sum_{a=1}^{A}(EU_a(opt) * w_a/(\sum_{o=1}^{O} EU_a(o)) \tag{1}$$

where $EU_d$(opt) is the utility value of an option ($opt$) at determinant $d$. $A$ is the set
of all ancestors of $d$ (i.e. determinants connect with $d$ at the previous level). $O$ is the
set of all available options. $w(a)$ is the weight of ancestor $a$. In this case, the weight
represents the importance of a decision-making determinant compare to others at
the same level and emphasizes on the heterogeneity of individuals. It also allows the
modeller to exclude determinants (i.e. setting their values to 0) that are not relevant
to a specific context. The combination process then continues until it reaches the
behaviour output list; the utility value of which can be translated to the probabilities
that an agent will choose that option. If the agent is assumed to be deterministic,
it picks the option that is correlated to the highest or lowest utility depending on
modeller's assumptions.

## 3.4  Summary of the Simulated Bounded Rationality Effects

With the two components above, we can summarise how the BRs can be simulated:

- **Sequential decision-making**: A decision-making cycle includes several steps,
  one after another. This procedure starts with the agent gathering information about
  the alternatives. Then, using its references, it filters, sorts, and cuts this list to a
  selected few. If triggered, these selected options are evaluated in the decision-
  making component. Finally, the highest/lowest evaluated alternative is selected
  and communicated to the environment. Using a procedural approach, this process
  follows the description of sequential decision-making in Sect. 1, i.e. the current
  step waits for the result of the previous step.
- **Emotional decision-making**: It is captured by the determinant *Affect* in the 2nd
  level of the *Decision-making* component (see Fig. 3). Its evaluation is dependent
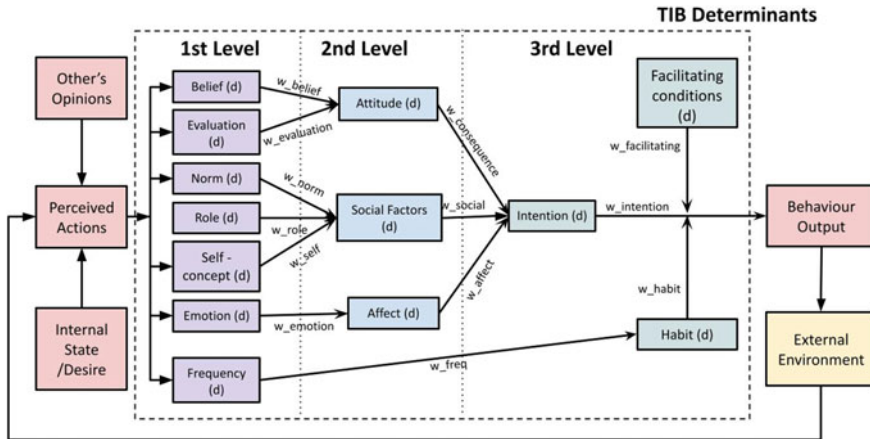
**Fig. 3** Decision-making component with TIB framework

on the context of decision-making. For example, our purchasing agent can rank
the level of comfort/pleasure it can have from a model compared to others. The
*Affect* determinant is associated with a weight ($w_{affect}$). By increasing this weight
and lowering the weights of other related determinants, we can highlight the con-
tribution of emotion to the overall behavioural output.

- **Habits**: Similar to emotion, the agent also accounts for past behaviour in its 3rd
  level of the TIB framework (see Fig. 3). Its weight can be adjusted to mark its
  influence on the final choice.
- **Multiple criteria**: The TIB framework in the *Decision-making* component allows
  users to capture different factors in decision-making, i.e. attitude(e.g. cost, time),
  norms, role, self-concept, emotion, habit, and past behaviour. A mapping with
  empirical data can be provided better to interpret these factors in a decision-
  making context. Function 1 provides a mean to combine them in the form of
  a utility value. Using associated weights, the agent can also decide which one has
  a larger/lower impact on the final choice. This concept also allows the agent to
  express its preferences on certain criteria of decision-making.
- **Confirmation bias**: In the *Perception* component, an agent filters the information
  received from the environment to form different short lists of options. This process
  represents the idea that the agent selects the information that supports its prefer-
  ence. The associated weights of each criterion mark the contribution of this bias
  to the final list. For example, in the car purchasing context, the user can generate
  an agent who only wants to receive information about electric cars by first setting
  the filter to only allow electric engine cars and zeroing all weights except for the
  engine's weight.
- **Bandwagon effect**: In its perception phase, the agent starts with observing its
  environment, including the patterns of its neighbour. It also accumulates the neigh-
  bours' opinions. This information is then used as a filter for in Perception compo-

nent (Fig. 2) and be fed into the Social factors determinant in the Decision-making component (Fig. 3). Each of them is associated with a weight to provide a way to compare its effects to other factors in the decision-making.

## 4 Case Study

This study focuses on observing the effect of bounded perception in an agent's decision-making. In the first Subsect. 4.1, we first calibrate our model using empirical data. The next subsection describes an experiment to demonstrate the function of the extended *Perception* component.

### 4.1 Data Mapping and Calibration

The environment in this study includes two main entities: *Market* and *Opinion Platform*. The *Market* consists of the details of the currently available car models, which are extracted from a Swiss car catalogue [20]. The given information include engine type, energy label, market price, brand and years of availability. The *Opinion Platform* provides reviews (value from 0 to 1) from the neighbourhood, dealer and media. Their weights are determined based on the network from the SHEDS panel data [28].

An agent in our model represents a household in Switzerland, which is generated using the process in [3]. There are currently 3080 agent profiles available. Each of them is associated with a weight to represent a portion of Switzerland's population. The behaviour outputs are multiplied by these weights to scale up to the national level.

In *Decision-making* component, the following properties can be mapped to the determinants of the first TIB level (see Fig. 3): Price—*Evaluation*, Review of dealer/media—*Role*, Review of neighbours—*Norm*, Brand of vehicle—*Self-concept*, Comfortability—*Emotion*, Availability of charging—*Facilitating condition* and Past usage of the same model—*Habit*.

To calibrate this purchasing model, two different sets of parameters corresponding to different components—*Perception* and *Decision-making*—are selected. In the Perception component, there are two main categories: thresholds for filters and weights (see Fig. 2). The thresholds include: (1) preferred engine (Gasoline, Diesel, Electric, Hybrid, other), (2) energy label (A, B, C, D, E, F and below), (3) price, (4) brand (1–8), recommendation level (value 0–1). In addition, each is associated with a weight, which also needs to be calibrated. In terms of the Decision-making component, we calibrate the following determinants' weights: price, energy label, recommendation, social status, brand, emotion, habit, attitude, social factor, intention and facilitating condition (charging infrastructure). At this stage of development, all weights will take a value in the set (0, 0.25, 0.5, 0.75, 1).

The number of parameters is significantly large, increasing the combined number of test runs exponentially. Therefore, we choose to perform a sensitive test for all parameters. The less critical parameters are assigned only two steps (0–1) in data calibration. From our tests, energy label, brand, and social status belong to this group. All parameters are then further grouped to create eight different agent purchasing profiles. Each agent is then assigned a random group for each of its parameters. This process ensures the heterogeneity in our agent population.
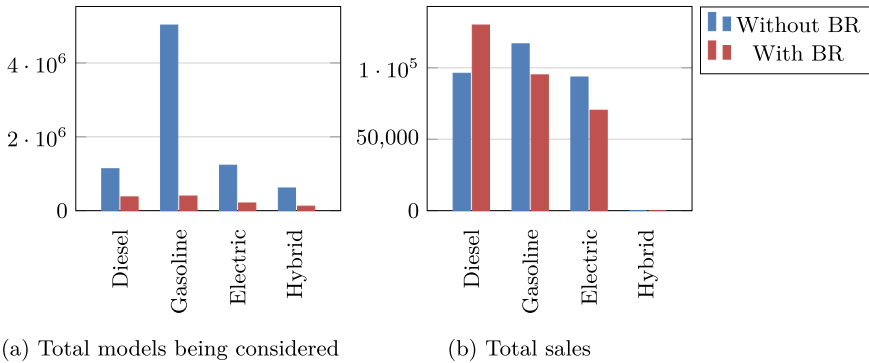
Our main objective is to minimise the error calculated by the total differences between the final number of vehicles purchased and real sales, multiplied by the weights (representing the adjusted importance) of the following criteria: (1) the total unit sales, (2) sum of sales of gasoline, diesel, electric and hybrid models and (3) the total sales of different clusters of models of different brands.

We calibrated with the data from 2015 to 2019. The more recent years, 2020–2021, are separated due to the effect of the pandemic COVID-19. Therefore, its car stock is adapted directly from correspondences in the SHEDS panel data. We repeat this procedure for all agent's profiles set at deterministic (i.e. choosing the best option) to find the smallest error. After a period of two weeks, the best setting satisfies the 1, 2, 3 condition with the yearly average errors after multiplied with weights equal to 305'485.

## 4.2   Evaluation of Bounded Perception

As the sequential, emotional, habitual, multi-criteria decision-making mechanisms are mainly implemented in the formerly developed *Decision-making* component, their effects on behaviour can be achieved by changing associated weights, similar to what was done previously in [16]. In this section, we focus on testing the functionality of the bounded *Perception* component in our agents. The number of vehicles (considered and purchased) calibrated for the final year (2019) is used as ground truth. We perform the experiment by turning the filtering, sorting and cutting functions off and evaluating the results against this ground truth. Figure 4a shows the results as the number of models being considered among the agent's population after the perception process. Figure 4b presents the final sales after the decision-making process. The figures are categorised by different engines, including diesel, gasoline, electric and hybrid vehicles.

The number of the models considered is much higher in the ground truth case (without BR), especially for gasoline models (considered nearly 12 times). When we apply filters with bounded perception, the distribution between different engine types is more balanced though it is proportioned to the case. In the total sales of ground truth, the highest number is for gasoline with 1.1 million vehicles. Even though electric vehicles are considered more, they have fewer sales. With the bounded perception applied, there are significant increases in the number of diesel cars sold. The gasoline and electric figures drop to 95'071 and 70'327 respectively. It is mainly due to better models of diesel and fewer models from the gasoline/electric type being

(a) Total models being considered

(b) Total sales

**Fig. 4** Simulation results in term of total number of vehicles per engine type

selected after the perception phase. Overall, we can clearly observe the difference in the number of models being considered (individual perception level) can lead to the difference in the percentage of car types sold (macro level).

## 5 Conclusion

In this study, we adopt our simulation platform—BedDeM—to simulate the impacts of different types of BRs. With the framework developed in [16, 17], sequential, emotional, habitual and multiple criteria decision-making can be considered in the agent's architecture. In addition, the Perception component is extended to cover the confirmation bias and bandwagon effect. This paper describes the agent's architecture design and provides an experiment to demonstrate the impact of bounded perception in the context of car purchasing in Switzerland. Similar experiments can be done to highlight the effect of single or combined BRs on an agent's decision-making and output.

The current model is still, however, missing some features, including variability in mapping between the first level determinants with SHEDS and MTMC data (see Fig. 3). This process can be accomplished by collaborating with a collaborator from the fields of economic or social science to derive a more accurate description of TIB's elements and generate more agent profiles to the current population.

There are also some promising research directions for our mobility platform. With the innovation in technology and increased environmental awareness, it has become more common for people to access electric or hydrogen vehicles. The model can provide a good indication of the roles of determinants in future scenarios (such as new infrastructures or government policies). Coupling with other models from different sectors can also provide a consumer's perspective where bounded rationalities can play a significant role in the agent's decision-making. As the topics provided in

Sect. 1 and their implementation in BedDeM are limited and simplified, one can implement more elaborate decision-making mechanisms in their modules to reflect the complexity of these topics.

# References

1. Andrighetto, G., Conte, R., Turrini, P., Paolucci, M.: Emergence in the loop: simulating the two way dynamics of norm innovation. In: Dagstuhl Seminar Proceedings. Schloss Dagstuhl-Leibniz-Zentrum für Informatik (2007)
2. Balke, T., Gilbert, N.: How do agents make decisions? A survey. J. Artif. Soc. Soc. Simul. **17**(4), 13 (2014)
3. Bektas, A., Schumann, R.: How to optimize gower distance weights for the k-medoids clustering algorithm to obtain mobility profiles of the swiss population. In: 2019 6th Swiss Conference on Data Science (SDS). pp. 51–56. IEEE (2019)
4. Boulouchos, K., Bach, C., Bauer, C., Bucher, D., Cerruti, D., Dehdarian, A., Filippini, M., Held, M., Hirschberg, S., Kannan, R., et al.: Pathways to a Net Zero co2 Swiss Mobility System: Sccer Mobility Whitepaper. Tech. rep, ETH Zurich (2021)
5. Carley, K.M., Gasser, L.: Computational organization theory. Multiagent systems: a modern approach to distributed artificial intelligence pp. 299–330 (1999)
6. Castro, J., Drews, S., Exadaktylos, F., Foramitti, J., Klein, F., Konc, T., Savin, I., van den Bergh, J.: A review of agent-based modeling of climate-energy policy. Wiley Interdisc. Rev. Clim. Change 11(4), e647 (2020)
7. Dignum, F., Dignum, V., Jonker, C.M.: Towards agents for policy making. In: International Workshop on Multi-Agent Systems and Agent-Based Simulation. pp. 141–153. Springer, Berlin (2008)
8. Frankish, K., Ramsey, W.M.: The Cambridge Handbook of Artificial Intelligence. Cambridge University Press (2014)
9. Georgeff, M., Pell, B., Pollack, M., Tambe, M., Wooldridge, M.: The belief-desire-intention model of agency. In: International Workshop on Agent Theories, Architectures, and Languages. pp. 1–10. Springer, Berlin (1998)
10. de Haan, P., Mueller, M.G., Scholz, R.W.: How much do incentives affect car purchase? Agent-based microsimulation of consumer choice of new cars-Part II: Forecasting effects of feebates based on energy-efficiency. Energy Policy **37**(3), 1083–1094 (2009)
11. Hahnel, U.J., Chatelain, G., Conte, B., Piana, V., Brosch, T.: Mental accounting mechanisms in energy decision-making and behaviour. Nat. Energy **5**(12), 952–958 (2020)
12. Jager, W., Janssen, M.: The need for and development of behaviourally realistic agents. In: International Workshop on Multi-Agent Systems and Agent-Based Simulation. pp. 36–49. Springer, Berlin (2002)
13. Kim, S., Lee, K., Cho, J.K., Kim, C.O.: Agent-based diffusion model for an automobile market with fuzzy topsis-based product adoption process. Expert Syst. Appl. **38**(6), 7270–7276 (2011)
14. Kiss, Á., Simonovits, G.: Identifying the bandwagon effect in two-round elections. Publ. Choice **160**(3), 327–344 (2014)
15. Kremmydas, D., Athanasiadis, I.N., Rozakis, S.: A review of agent based modeling for agricultural policy evaluation. Agric. Syst. **164**, 95–106 (2018)
16. Nguyen, K., Schumann, R.: An exploratory comparison of behavioural determinants in mobility modal choices. In: Conference of the European Social Simulation Association. pp. 569–581. Springer, Berlin (2019)

17. Nguyen, K., Schumann, R.: On developing a more comprehensive decision-making architecture for empirical social research: Agent-based simulation of mobility demands in Switzerland. In: International Workshop on Multi-Agent Systems and Agent-Based Simulation. pp. 39–54. Springer, Berlin (2019)
18. Nickerson, R.S.: Confirmation bias: a ubiquitous phenomenon in many guises. Rev. Gen. Psychol. **2**(2), 175–220 (1998)
19. Pereira, D., Oliveira, E., Moreira, N., Sarmento, L.: Towards an architecture for emotional bdi agents. In: 2005 Portuguese Conference on Artificial Intelligence. pp. 40–46. IEEE (2005)
20. Revue Automobile Catalogue. https://revueautomobile.ch/, [Last accessed 28 Apr 2022]
21. Rollwage, M., Fleming, S.M.: Confirmation bias is adaptive when coupled with efficient metacognition. Philos. Trans. Royal Soc. B **376**(1822), 20200131 (2021)
22. Russell, S., Norvig, P.: Artificial Intelligence: A Modern Approach, 3rd edn. Prentice Hall Press, USA (2010)
23. Simon, H.A.: Rationality as process and as product of thought. Am. Econ. Rev. **68**(2), 1–16 (1978)
24. Simon, H.A.: Making management decisions: the role of intuition and emotion. Acad. Manag. Perspect. **1**(1), 57–64 (1987)
25. Sobkowicz, P.: Opinion dynamics model based on cognitive biases of complex agents. J. Artif. Soc. Soc. Simul. 21(4) (2018)
26. Thiriot, S., Kant, J.D.: A multi-agent cognitive framework to model human decision making under bounded rationality. In: IAREPSABE 06 (International Conference on Behavioural Economics and Economic Psychology), Paris (2006)
27. Urban, C.: Pecs: A reference model for the simulation of multi-agent systems. In: Tools and Techniques for Social Science Simulation, pp. 83–114. Springer, Berlin (2000)
28. Weber, S., Burger, P., Farsi, M., Martinez-Cruz, A.L., Puntiroli, M., Schubert, I., Volland, B.: Swiss Household Energy Demand Survey (SHEDS): Objectives, Design, and Implementation. Tech. rep., IRENE Working Paper (2017)

# The Benefits of Coordination in Adaptive Virtual Teams

**Darío Blanco-Fernández, Stephan Leitner, and Alexandra Rausch**

**Abstract** The emergence of new organizational forms—such as virtual teams—has brought forward some challenges for teams. One of the most relevant challenges is coordinating the decisions of team members who work from different places. Intuition suggests that task performance should improve if the team members' decisions are coordinated. However, previous research suggests that the effect of coordination on task performance is ambiguous. Specifically, the effect of coordination on task performance depends on aspects such as the team members' learning and the changes in team composition over time. This paper aims to understand how these two factors moderate the relationship between coordination and task performance. We implement an agent-based modeling approach based on the *NK* framework to fulfill our research objective. Our results suggest that both factors have moderating effects. Specifically, we find that excessive individual learning harms the task performance of fully autonomous teams but is less detrimental for teams that coordinate their decisions. In addition, we find that teams that coordinate their decisions benefit from changing their composition in the short term, but fully autonomous teams do not. In conclusion, teams that coordinate their decisions benefit more from individual learning and dynamic composition than teams that do not coordinate. Nevertheless, we should note that the existence of moderating effects does not imply that coordination improves task performance. Whether coordination improves task performance depends on the interdependencies between the team members' decisions.

**Keywords** Coordination · Complex task · Individual learning · Team composition · Agent-based modeling

D. Blanco-Fernández (✉) · S. Leitner · A. Rausch
University of Klagenfurt, Klagenfurt 9020, Austria
e-mail: dario.blanco@aau.at

S. Leitner
e-mail: stephan.leitner@aau.at

A. Rausch
e-mail: alexandra.rausch@aau.at

435

# 1 Introduction

The COVID-19 pandemic and the development of communication technologies have greatly expanded teleworking arrangements in organizations. New forms of collaboration—such as virtual teams—have emerged in response to these changes in the workplace [12, 15, 18]. Virtual teams are "*groups of people with a common purpose who carry out interdependent tasks across locations and time, using technology to communicate much more than they use face-to-face meetings*" [5].

Virtual teams often rely on self-organization and decentralized decision-making to achieve a shared or centralized objective [8, 12, 15, 18]. For example, an organization might try to improve the functioning of a particular software by giving their employees autonomy in team formation and decision-making [12]. The success of virtual teams is usually determined by their members' initiative in searching for and implementing new solutions to the task they face [15]. Thus, virtual team members might improve task performance by *learning* about the task and adapting their knowledge to its requirements [4, 17].

The goal of virtual teams is often to complete a particular task, i.e., virtual teams are often task-oriented [15, 18]. This task-oriented perspective implies that virtual teams often exhibit a *dynamic composition*, i.e., their composition changes over time [15]. As the task is being completed, new challenges and demands might appear. Virtual teams often change their composition to adapt to these challenges [1, 15, 19]. By changing their composition, virtual teams aim to integrate new members that bring knowledge previously unavailable to the team [3, 4, 17].

The decentralized nature of virtual teams and their dynamic composition implies that virtual team members might struggle to communicate their intended actions effectively to other team members [5, 8, 12, 15]. Prior research suggests that virtual team members should coordinate their actions to avoid this problem [15]. Coordination might be achieved by establishing clear rules and procedures which ensure communication and regulate individual decision-making [15, 18].

Although intuition suggests that coordinating the team's decisions should increase task performance [7], prior results indicate that the effect of coordination on task performance is not straightforward [16]. Coordination does not unfold positive effects if tasks are simple but is beneficial for sufficiently complex tasks [16, 21]. In addition, coordination allows teams to take full advantage of their members' learning, so it is particularly beneficial when team members have access to a large set of solutions to the task [6, 7, 16]. However, prior research usually focuses on teams which do not change their composition [16, 21]. Consequently, previous research ignores the moderating effects of dynamic team composition in the relationship between coordination and task performance. Our paper aims to fill this research gap.

This paper focuses on self-organized teams—such as virtual teams—that solve complex tasks. Our objective is to understand how the effect of coordination on task performance is moderated by (i) individual learning and (ii) team composition. To achieve this research objective, we build on previous research by introducing coordination between the team members' decisions [2–4]. We contribute to the literature
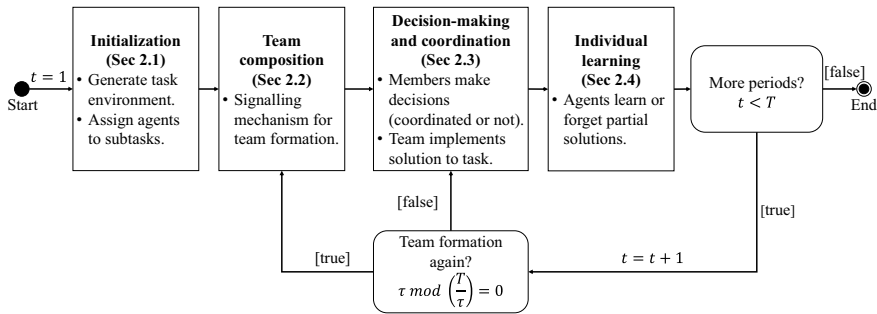
**Fig. 1** Temporal sequence of the model

by showing how coordination does not only allow teams to fully grasp the benefits of individual learning but also those of dynamic composition. Additionally, our results also provide practical advice for virtual teams regarding coordination. The remainder of this paper is organized as follows: Sect. 2 provides details on the model. Section 3 provides a description and discussion of the main results. Finally, Sect. 4 concludes the paper.

## 2 The Model

We implement an agent-based model based on the *NK* framework [10].[1] The model consists of four building blocks: The task environment and the agents (see Sect. 2.1), team formation (see Sect. 2.2), decision-making and coordination (see Sect. 2.3), and individual learning (see Sect. 2.4). The four building blocks correspond to the sequence of events of the model, which we illustrate in Fig. 1.

### 2.1 Initialization

**Task environment** The complex task that teams face consists of $N = 12$ binary interdependent decisions. We divide the $N$-dimensional complex task into $M = 3$ subtasks of equal length $S = N/M = 4$. We denote each subtask by a vector $\mathbf{d}_m = (d_{S \cdot (m-1)+1}, \ldots, d_{S \cdot m})$ and the complex task by the vector $\mathbf{d} = (d_1, \ldots, d_N)$.[2]

Each decision $d_n \in [0, 1]$ contributes $c_n$ to task performance $C(\mathbf{d})$. This contribution depends on the decision itself and $K$ other decisions, so $c_n = f(d_n, d_{i_1}, \ldots, d_{i_K})$, where $\{i_1, \ldots, i_K\} \subseteq \{1, \ldots, n-1, n+1, \ldots, N\}$ and $0 \leq K \leq N - 1$. We ran-

---

[1] The model has been implemented in Python 3.7.4.

[2] It follows that $\mathbf{d}_1 \frown \cdots \frown \mathbf{d}_M = \mathbf{d}$, where $\frown$ is the concatenation of each subtask.

domly generate contributions using a uniform distribution, $c_n \sim U(0, 1)$. The overall task performance is the average of all contributions $C(\mathbf{d}) = \frac{1}{N} \sum_{n=1}^{N} c_n$.

Each possible vector of $N = 12$ binary values is a *solution* to the complex task and has an associated performance. There are $2^S = 16$ possible *partial solutions* to each subtask and $2^N = 4.096$ possible solutions to the complex task. The mapping of each solution to its associated performance is the *performance landscape*. The team moves gradually on the performance landscape, following a steepest ascent hill-climbing search for new, better-performing solutions.

The *task complexity*—determined by $K$—partly influences the success of the search process. The higher $K$, the more complex the task is, and the more rugged the performance landscape [10]. Several local maxima characterize a rugged performance landscape. Consequently, the higher the task complexity, the more likely it is for teams to get stuck at suboptimal solutions [10]. Regarding task complexity, we consider two different scenarios: Low ($K = 3$) and moderate complexity ($K = 5$).

The *interdependence pattern* also affects the performance landscape's shape [13]. The interdependence pattern reflects which contributions depend on which decisions. We consider three interdependence patterns, which we represent in Fig. 2:
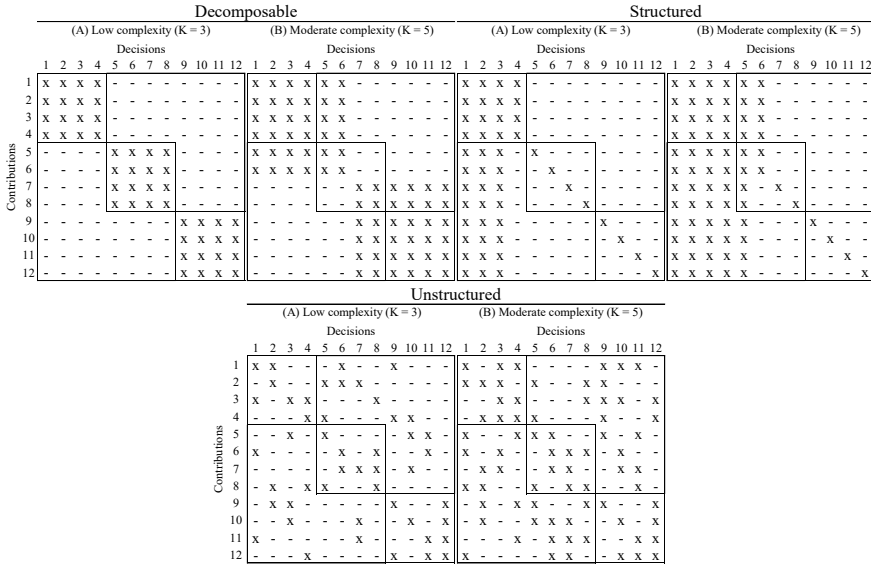
- *Decomposable*: Interdependencies are shaped in squares of size $K + 1$. For $K = 3$, the task is *perfectly decomposable*, and there are no interdependencies between subtasks. For $K = 5$, by contrast, there are interdependencies between subtasks.
- *Structured*: The $K$ first decisions affect the remaining contributions. Thus, there is one subtask that heavily influences task performance.
- *Unstructured*: Interdependencies between decisions are randomly allocated throughout the task, resulting in interdependencies between subtasks in all cases.

**Agents** To solve the complex task, a team is formed by choosing one member per subtask—i.e., $M = 3$ members—out of a population of $P = 30$ agents. These agents are heterogeneous and have *limited capabilities* concerning the complex task. We limit the agents' capabilities in two ways. First, each agent can only solve one subtask $m$. Second, in the first period, we endow each agent with just one random partial solution to subtask $m$. Agents must solve the entire complex task to experience positive utility.

Agents are myopic as they aim to optimize just their immediate utility. Only team members experience positive utility.[3] The utility function of an agent assigned to subtask $m$ is the weighted sum of their own performance contributions $C(\mathbf{d}_{mt})$ and the performance contributions of the *residual decisions* $C(\mathbf{d}_{rt})$, where $r = \{1, \ldots, M\} \in \mathbb{N}$ and $r \neq m$. We denote the residual decisions by $\mathbf{D}_{mt} = (\mathbf{d}_{1t}, \ldots, \mathbf{d}_{\{m-1\}t}, \mathbf{d}_{\{m+1\}t}, \ldots, \mathbf{d}_{Mt})$. Agent $m$'s utility is calculated using Eq. 1:

$$U(\mathbf{d}_{mt}, \mathbf{D}_{mt}) = \frac{1}{2} \cdot \left( C(\mathbf{d}_{mt}) + \cdot \frac{1}{M-1} \sum_{\substack{r=1 \\ r \neq m}}^{M} C(\mathbf{d}_{rt}) \right). \tag{1}$$

---

[3] Agents who do not join the team in one period get utility equal to 0.

**Fig. 2** Interdependence matrices. Each matrix depicts which contributions (*y*-axes) depend on which decisions (*x*-axes). Interdependencies are indicated with an *X*. Since a contribution depends on its own decision, there is an *X* in each element of the main diagonal. Solid lines indicate the subtasks

## 2.2 Team Composition

Agents always have an incentive to participate in the team since it is the only way to experience positive utility. We assume that agents are fully aware of how team formation works and that they do not cheat. Additionally, we omit communication between agents during team formation. These assumptions assure that agents do not behave strategically or form beliefs about other agents. Finally, we assume that one agent per subtask is sufficient to solve the complex task. Consequently, three members form the team.

The objective of team formation is to assure that the current team members are the best-available agents for solving the task at any given period. Team formation works as follows. The set of solutions agent $m$ knows is $\mathbf{S}_m = \left( \hat{\mathbf{d}}_{m1}, \ldots, \hat{\mathbf{d}}_{mI} \right)$ where $\hat{\mathbf{d}}_{mi}$ is a solution to subtask $\mathbf{d}_m$, $i = \{1, \ldots, I\} \in \mathbb{N}$ and $1 \leq I \leq 2^S$. Each agent estimates the utility for each solution they know, i.e., $\forall \hat{\mathbf{d}}_{mi} \in \mathbf{S}_{mt}$. Since we omit communication, agents use the residual decisions from the previous period $\mathbf{D}_{m\{t-1\}}$ as a basis for their estimations. Agent $m$'s *estimated utility* is then:

$$\mathrm{EU}(\mathbf{d}_{mt}, \mathbf{D}_{m\{t-1\}}) = \frac{1}{2} \cdot \left( C(\mathbf{d}_{mt}) + \cdot \frac{1}{M-1} \sum_{\substack{r=1 \\ r \neq m}}^{M} C(\mathbf{d}_{r\{t-1\}}) \right) + e; \qquad (2)$$

where $e$ is an error term which follows a normal distribution $e \sim N(0, 0.01)$. This error term reflects the mistakes that team members might make when estimating the effects of their decisions [9, 22].

After the estimation, each agent signals the highest estimated utility $U(\hat{\mathbf{d}}^*_{mt}, \mathbf{D}_{m\{t-1\}})$, where $\hat{\mathbf{d}}^*_{mt} := \arg\max_{\mathbf{d}' \in \mathbf{S}_{mt}} U(\mathbf{d}', \mathbf{D}_{m\{t-1\}})$ is the solution that maximizes agent $m$'s estimated utility at time $t$. The agent who signals the highest estimated utility for each subtask $m$ becomes a team member.

The agents form the first team iteration in the first period. Afterwards, team formation is repeated every $\tau$ periods. The higher (lower) $\tau$, the less (more) frequently a team changes its composition. We study three different scenarios for $\tau$:

- Teams with a long-term composition do not change their composition over time. We denote this scenario by $\tau = \emptyset$.
- Teams with a medium-term composition change their composition every $\tau = 10$ periods.
- Teams with a short-term composition change their composition at every period, i.e., $\tau = 1$.

## 2.3   Decision-Making and Coordination

The three team members choose a team solution to the complex task at every period. In the benchmark scenario, members of a *fully autonomous* team make their choices independently and simultaneously [16, 21]. They calculate the estimated utility for every partial solution they know, i.e. $\forall \hat{\mathbf{d}}_{mi} \in \mathbf{S}_{mt}$, following Eq. 2. Each member's choice is $\hat{\mathbf{d}}^*_{mt}$, i.e., the partial solution associated with the highest estimated utility. Finally, the concatenation of all member's choices is the team solution for the current period $\mathbf{d}_t := \hat{\mathbf{d}}^*_{1t} \frown \cdots \frown \hat{\mathbf{d}}^*_{Mt}$.

We contrast this benchmark scenario with a scenario in which the team members coordinate their decisions. The coordination mechanism is based on the *liaison* organizational archetype described in [16]. We assume that all agents know how the coordination mechanism works. Initially, each team member ranks all partial solutions they know $\hat{\mathbf{d}}_{mi} \in \mathbf{S}_{mt}$ regarding their estimated utility (see Eq. 2). Then, each member chooses the two highest partial solutions $\hat{\mathbf{d}}^{(1)}_{mt}$ and $\hat{\mathbf{d}}^{(2)}_{mt}$, where the solution with the highest expected utility is ranked first and the solution with the second highest expected utility is ranked second. The team members bring these partial solutions to a coordination session.

Two candidate solutions are constructed in order, first by concatenating the preferred choices and then the second-preferred choices so $\mathbf{d}^{(j)}_t := \hat{\mathbf{d}}^{(j)}_{1t} \frown \cdots \frown \hat{\mathbf{d}}^{(j)}_{Mt}$ where $\hat{\mathbf{d}}^{(j)}_{mt}$ is agent $m$'s $j$th preferred choice. Each team member sequentially evaluates the two candidate solutions $\mathbf{d}^{(j)}_t$ regarding their estimated utility. If the estimated utility of a candidate solution is higher than the last achieved utility, the team member accepts the candidate solution, i.e., they accept the solution if $EU_m(\mathbf{d}^{(j)}_t) > U_m(\mathbf{d}_{t-1})$. Otherwise, they veto it. The veto from one member is enough to reject the candidate

solution. If all members accept a candidate solution, it is chosen as the team solution for the current period, so $\mathbf{d}_t = \mathbf{d}_t^{(j)}$. Conversely, the team solution remains constant from the previous period if members veto both candidate solutions, so $\mathbf{d}_t = \mathbf{d}_{t-1}$.

### 2.4  Individual Learning

Agents overcome their limited capabilities by *learning* about subtask $m$. Specifically, agents learn by exploring the solution space and changing their set of partial solutions $\mathbf{S}_m$ over time [10]. Learning occurs at the end of each period for all agents—even if they are not part of the group—and consists of two separate mechanisms. First, with probability $\mathbb{P}$, agents might *discover* a partial solution. The partial solution they discover differs in the value of one decision from any currently-known partial solution. For example, if agent $m$'s solution space $\mathbf{S}_{mt}$ consists only of decision $\hat{\mathbf{d}}_{mt} = (0, 0, 0, 0)$, then agent $m$ can learn one of these four solutions at time $t$: $(0, 0, 0, 1)$, $(0, 0, 1, 0)$, $(0, 1, 0, 0)$, and $(1, 0, 0, 0)$. Second, with the same probability $\mathbb{P}$, agents might *forget* a partial solution that is not utility-maximizing in the current period, i.e., they can forget any solution in $\mathbf{S}_{mt}$ except $\hat{\mathbf{d}}_{mt}^*$.[4] We study probabilities of learning between $\mathbb{P} = 0$ and $\mathbb{P} = 1$ in intervals of 0.1.

### 2.5  Parameters and Performance Measures

Our research comprises 396 different scenarios. Each scenario consists of 1500 simulation rounds of $T = 200$ periods each. We summarize the main parameters of the model and their values in Table 1.

 We normalize the observed task performance at each period $C(\mathbf{d_t})$ by the maximum achievable performance at each simulation round $C^*$. Normalization assures that we can compare different scenarios in terms of task performance.

 We train regression tree models using the normalized task performance as the dependent variable and the independent variables of Table 1. We then compute partial dependencies between task performance and the moderating factors, i.e., individual learning and team composition. To calculate partial dependencies, we first define $\mathbf{X}$ as the set of all independent variables. The set $\mathbf{X}$ is divided into two subsets. Subset $\mathbf{X}^s$ corresponds to the scope variable, i.e., individual learning or dynamic team composition. Subset $\mathbf{X}^c$ includes the remaining independent variables.[5] We compute the partial dependence of task performance on the moderating factor studied according to $f^s(\mathbf{X}^s) = E_c(f(\mathbf{X}^s, \mathbf{X}^c)) \approx \frac{1}{V} \sum_{i=1}^{V} f(\mathbf{X}^s, \mathbf{X}_{(i)}^c)$, where $V$ is the number of independent variables in $\mathbf{X}^c$ and $\mathbf{X}_{(i)}^c$ corresponds to each variable. We employ this method to understand the patterns related to our research objective [11].

---

[4] An agent which only knows one solution cannot forget it.

[5] It follows that $\mathbf{X}^s \cup \mathbf{X}^c = \mathbf{X}$.

**Table 1** Parameters

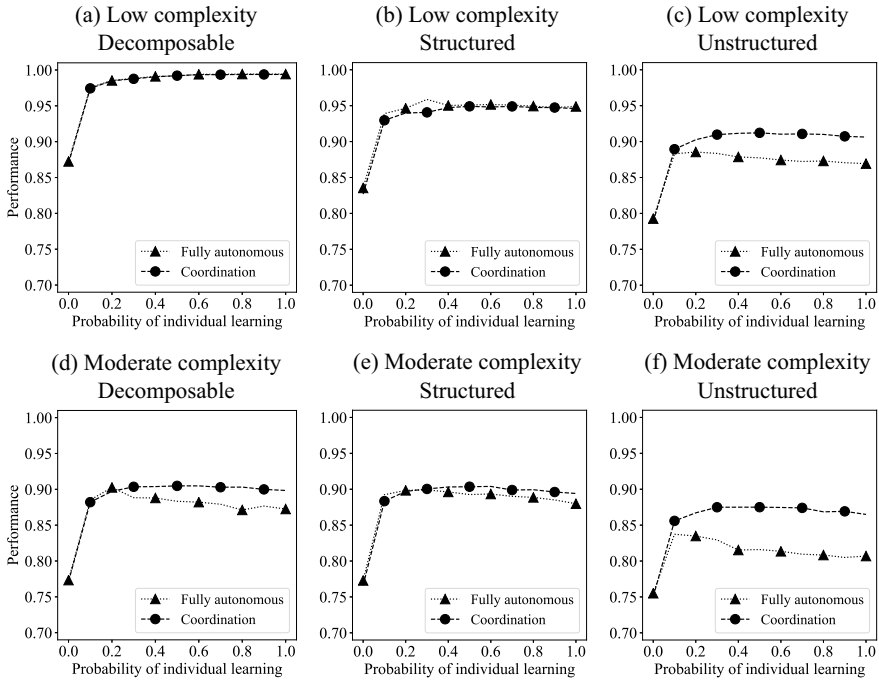| Type | Variables | Notation | Values |
|---|---|---|---|
| Independent variables | Task complexity | $K$ | {3, 5} |
| | Interdependence structure | *Matrix* | See Fig. 2 |
| | Team composition | $\tau$ | {Ø, 1, 10} |
| | Learning probability | $\mathbb{P}$ | 0:0.1:1 |
| | Time period | $t$ | 1:1:100 |
| | Coordination | N/A | *Fully autonomous, coordination* |
| Dependent variable | Task performance | $C(\mathbf{d_t})$ | [0, 1] |
| Other parameters | Number of decisions | $N$ | 12 |
| | Population of agents | $P$ | 30 |
| | Number of subtasks | $M$ | 3 |
| | Number of simulations | $\Phi$ | 1500 |
| | Error term | $e$ | $e \sim N(0, 0.01)$ |

## 3 Results and Discussion

According to prior research, coordination is more beneficial for teams that search extensively for new solutions [14, 16]. Teams may acquire new solutions because (i) their members learn about the complex task or (ii) they change their composition [17]. Some authors advocate that research on coordination should consider search processes at multiple levels [20]. Our research follows this suggestion and aims to understand how individual learning ($\mathbb{P}$) and dynamic team composition ($\tau$) moderate the relationship between coordination and task performance. In the following sub-sections, we separately study the moderating effects of individual learning (Sect. 3.1) and team composition (Sect. 3.2).

### 3.1 Moderating Effect of Individual Learning

To study the moderating effect of individual learning, we calculate partial dependencies using the individual learning probability as the scope variable. We represent the results in Fig. 3. Each plot shows the partial dependence of task performance on the learning probability for each level of complexity $K$ and for each interdependence structure.

The moderating effect of individual learning depends on the level of complexity and the interdependence structure studied. There is no moderating effect of individual learning for decomposable and structured tasks of low complexity (Fig. 3a, b). The effect of increasing the individual learning probability on task performance is remarkably similar for teams that coordinate their decisions and fully autonomous
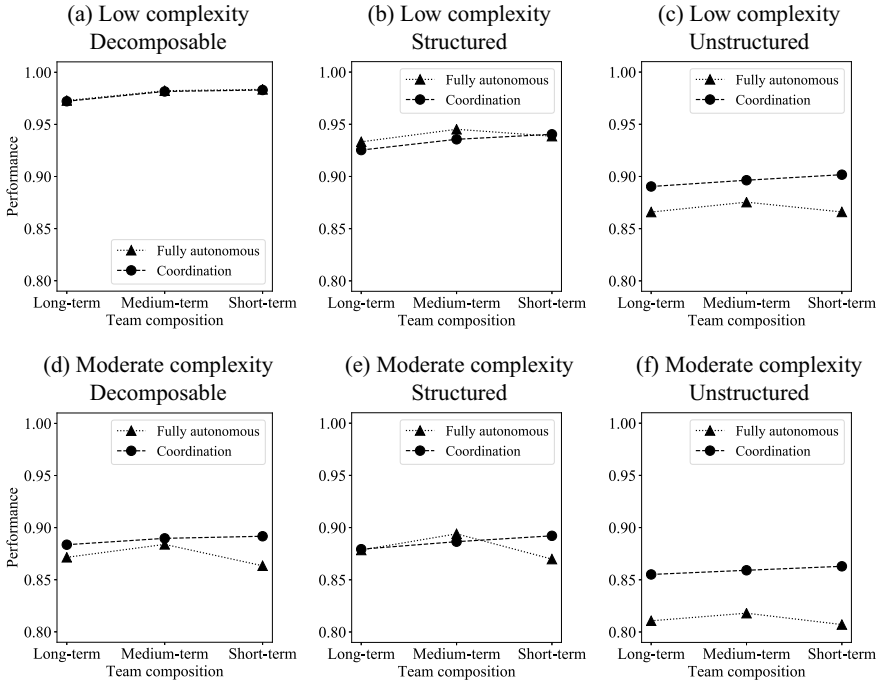
**Fig. 3** Partial dependencies of task performance ($y$-axes) on the probability of individual learning ($x$-axes)

teams. Task performance reacts strongly to initial increases in the probability of individual learning. However, the task performance stabilizes for higher values of $\mathbb{P}$.

Individual learning has a moderating effect for unstructured tasks of low complexity (Fig. 3c) and moderately complex tasks (Fig. 3d–f). At relatively low levels of learning, the effect of increasing $\mathbb{P}$ is similar for teams that coordinate their decisions and fully autonomous teams. Task performance increases, albeit the marginal positive effect decreases with each increase in $\mathbb{P}$. Eventually, there is a threshold value for the individual learning probability at which its effect on task performance turns negative. This negative effect is highly relevant for fully autonomous teams. In contrast, the decrease in task performance is barely notable in teams that coordinate their decisions. Consequently, the benefits of coordination increase with individual learning.

Our results show that agents learning more does not necessarily increase task performance, but might even harm it (Fig. 3c–f). Therefore, our results align with previous research, which highlights the negative impact that excessive individual learning might have on task performance [2–4]. In addition, we show that coordination reduces these adverse effects notably. The reason for this might lie in the evaluation of the team members' solutions. Coordination allows teams to evaluate their members' partial solutions more efficiently, increasing task performance [7].

**Fig. 4** Partial dependencies of task performance ($y$-axes) on team composition ($x$-axes)

While individual learning might be a key determinant of the performance of virtual teams, the team members should coordinate to grasp the full benefits of individual learning [6, 7]. This is particularly relevant for virtual teams, as their performance depends on their members' ability to find and implement new solutions [15].

## 3.2  Moderating Effect of Team Composition

We compute partial dependencies using team composition as the scope variable to study the second moderating effect. Each plot in Fig. 4 shows the partial dependence of task performance on team composition for each level of complexity $K$ and each interdependence structure considered.

In perfectly decomposable tasks, we find that team composition $\tau$ has the same impact on task performance regardless of coordination (see Fig. 4a). A dynamic team composition slightly improves task performance compared to a stable team composition, irrespective of the frequency of team formation.

For the remaining scenarios (Fig. 4b–f), there are different patterns depending on coordination. The task performance of teams that coordinate their decisions increases with the frequency of team formation. Consequently, teams that coordinate their

decisions benefit more from changing their composition in the short term. By contrast, fully autonomous teams only benefit from changing their composition in the medium term. Changes in the short term might unfold neutral (see Fig. 4b, c, f) or negative effects (see Fig. 4d, e) on task performance.

Teams that change their composition find more solutions than stable teams [2], but they require coordination to translate this search process into improvements in task performance [16, 19]. The lack of coordination between the team members' decisions makes fully autonomous teams unable to benefit from changing their composition in the short term. In contrast, teams that coordinate their decisions benefit from more frequent changes in their composition (see Fig. 4b–f). Thus, our results align with the insights from prior research [16, 19].

Whether this moderating effect results in coordination improving task performance depends on the interdependence structure. Coordination never improves task performance for structured tasks of low complexity (Fig. 4b). For decomposable and structured tasks of moderate complexity (Fig. 4d, e), coordination only improves task performance if teams have a short-term composition. Finally, coordination is always beneficial for unstructured tasks (Fig. 4c, f). Additionally, the positive effect of coordination in unstructured tasks is higher for short-term team composition.

Our results align partially with the insights given by Tannenbaum et al. [19], who claim that dynamic teams should coordinate their decisions to improve task performance. We show that coordination is beneficial only in certain situations. The benefits of coordination grow with the frequency of changing composition and the interdependencies between the team members' decisions. According to prior research, teams might achieve indirect coordination by assigning the decisions in such a manner that the interdependencies between subtasks are minimized [14]. If subtasks are less interdependent, there is less need for coordination, as the effects of the members' decisions on the remaining members' contributions diminish [21]. This indirect coordination is reflected in decomposable and structured tasks. Consequently, the impact on task performance of coordinating the team decisions is reduced [14, 21].

Coordination might be advisable for virtual teams, as their composition is generally more dynamic than in traditional teams [15]. We should note, however, that coordination might not be necessary if tasks are relatively simple. In this case, coordination might only slow down decision-making [16], which is particularly relevant for task-oriented groups—such as virtual teams [15]. Thus, virtual team members should carefully consider how interdependent their actions are before coordinating.

## 4 Concluding Remarks

This paper studies how individual learning and team composition moderate the relationship between coordination and task performance. We find moderating effects of both variables that are more relevant, the more interdependent the team members' decisions. Regarding individual learning, our results suggest that the adverse effects of increasing individual learning excessively are reduced if teams coordinate their

decisions. Regarding team composition, we find that fully autonomous teams only benefit from changing their composition in the medium term. In contrast, teams that coordinate their decisions benefit more from changing their composition in the short term.

Coordination does not affect the performance of decomposable and structured tasks of low complexity. In contrast, coordination improves task performance for (i) unstructured tasks, (ii) for sufficiently high individual learning, or (iii) for a short-term composition. From a practical standpoint, our results suggest that coordination might benefit virtual teams, as they rely on individual learning to improve task performance, and their composition is often dynamic. Virtual team members, however, should carefully consider the interdependencies between their decisions before implementing coordination.

Our research has some limitations. We only test the effect of one coordination mechanism. More coordination mechanisms could be added to extend this research, as in [16]. Additionally, prior research suggests that individual learning and team composition interact [4]. Future research can study the joint moderating effect of individual learning and team composition on the relationship between coordination and task performance. Finally, researchers have suggested other aspects that might affect the relationship between team composition and coordination. In particular, researchers cite team identity as a relevant factor for coordination in dynamic teams [15, 18]. Future extensions of our research could consider this aspect.

# References

1. Bell, S.T., Outland, N.: Team composition over time. Res. Managing Groups Teams **18**, 3–27 (2017)
2. Blanco-Fernández, D., Leitner, S., Rausch, A.: Autonomous group formation of heterogeneous agents in complex task environments. In: Czupryna, M., Kamiński, B. (eds.) Advances in Social Simulation, pp. 131–144. Springer International Publishing, Cham (2022)
3. Blanco-Fernández, D., Leitner, S., Rausch, A.: Multi-level adaptation of distributed decision-making agents in complex task environments. In: Van Dam, K.H., Verstaevel, N. (eds.) Multi-agent-Based Simulation XXII, pp. 29–41. Springer International Publishing, Cham (2022)
4. Blanco-Fernández, D., Leitner, S., Rausch, A.: Dynamic groups in complex task environments: to change or not to change a winning team? (2022). https://arxiv.org/abs/2203.09157
5. Cramton, C.D.: The mutual knowledge problem and its consequences for dispersed collaboration. Organ. Sci. **12**(3), 346–371 (2001)
6. Edmondson, A.C., Bohmer, R.M., Pisano, G.P.: Disrupted routines: team learning and new technology implementation in hospitals. Adm. Sci. Q. **46**(4), 685–716 (2001)
7. Edmondson, A.C., Dillon, J., Roloff, K.: Three perspectives on team learning: outcome improvement, task mastery, and group process. Acad. Manage. Ann. **1**(1), 269–314 (2007)
8. Foster, I., Kesselman, C., Tuecke, S.: The anatomy of the grid: enabling scalable virtual organizations. Int. J. High Perform. Comput. Appl. **15**(3), 200–222 (2001)
9. Leitner, S., Wall, F.: Simulation-based research in management accounting and control: an illustrative overview. J. Manage. Control **26**(2–3), 105–129 (2015)
10. Levinthal, D.A.: Adaptation on rugged landscapes. Manage. Sci. **43**(7), 934–950 (1997)

11. Patel, M.H., Abbasi, M.A., Saeed, M., Alam, S.J.: A scheme to analyze agent-based social simulations using exploratory data mining techniques. Complex Adapt. Syst. Model. **6**(1), 1–17 (2018)
12. Puranam, P., Alexy, O., Reitzig, M.: What's "new" about new forms of organizing? Acad. Manage. Rev. **39**(2), 162–180 (2014)
13. Rivkin, J.W., Siggelkow, N.: Patterned interactions in complex systems: implications for exploration. Manage. Sci. **53**(7), 1068–1085 (2007)
14. Rivkin, J.W., Siggelkow, N.: Balancing search and stability: interdependencies among elements of organizational design. Manage. Sci. **49**(3), 290–311 (2003)
15. Saunders, C.S., Ahuja, M.K.: Are all distributed teams the same? Differentiating between temporary and ongoing distributed teams. Small Group Res. **37**(6), 662–700 (2006)
16. Siggelkow, N., Rivkin, J.W.: Speed and search: designing organizations for turbulence and complexity. Organ. Sci. **16**(2), 101–122 (2005)
17. Simon, H.A.: Bounded rationality and organizational learning. Organ. Sci. **2**(1), 125–134 (1991)
18. Squicciarini, A.C., Paci, F., Bertino, E.: Trust establishment in the formation of virtual organizations. Comput. Stan. Interfaces **33**(1), 13–23 (2011)
19. Tannenbaum, S.I., Mathieu, J.E., Salas, E., Cohen, D.: Teams are changing: are research and practice evolving fast enough? Ind. Organ. Psychol. **5**(1), 2–24 (2012)
20. Tošić, P.T., Vilalta, R.: A unified framework for reinforcement learning, co-learning and meta-learning how to coordinate in collaborative multi-agent systems. Procedia Comput. Sci. **1**(1), 2217–2226 (2010)
21. Wall, F.: Emergence of task formation in organizations: balancing units' competence and capacity. J. Artif. Soc. Soc. Simul. **21**(2), 1–26 (2018)
22. Wall, F., Leitner, S.: Agent-based computational economics in management accounting research: opportunities and difficulties. J. Manage. Acc. Res. (11) (2021)

# The Equity Premium Puzzle: An Application of an Agent-Based Evolutionary Model

**Luca Gerotto, Paolo Pellizzari, and Marco Tolotti**

**Abstract** We describe an agent-based model of a financial market with a stock and a bond. Agents compete in repeated rounds, decide whether to acquire costly information and can pick one of 16 strategies to allocate their investments, under evolutionary pressure driven by the comparison of the realized short-term revenues from trading. We show that, while informed traders survive in some cases, the equilibrium shares are strongly biased in favor of strategies that make little use of information and systematically overestimate the riskiness of the stock. As a consequence, the majority of the population ends up in buying fewer stocks than would be otherwise expected or deemed rational. This evolutionary dynamics offers a novel way to explain the equity premium puzzle first described by Mehra and Prescott (*The equity premium: A puzzle*. Journal of Monetary Economics 1985), according to which it's hard to find reasons for the widespread lack of investment in risky assets. Evolution based on a straightforward comparison of revenues is a simple and cognitively appealing avenue to reach a population of traders using (over-)cautious strategies to curb the risk of long-term "financial extinction". Simulations run in NetLogo also demonstrate that very little information may be used in noisy markets or when the cost of information is substantial.

**Keywords** Agent-based models · Evolutionary game theory · Equity premium puzzle

L. Gerotto (✉)
Faculty of Economics, Department of Economics and Finance, Universitá Cattolica del Sacro Cuore, Milan, Italy
e-mail: luca.gerotto@unicatt.it
URL: https://sites.google.com/view/gerottoluca/home-page

P. Pellizzari
Department of Economics, Ca' Foscari University of Venice, Venice, Italy

M. Tolotti
Department of Management, Ca' Foscari University of Venice, Venice, Italy

449

# 1 Introduction

As famously pointed out in Mehra and Prescott [6] it is difficult to reconcile standard financial economic models with the observation that investors purchase relatively small amounts of stocks, whose average returns are historically much bigger than the safe rate (obtained with highly-rated bonds or bills). This conundrum has been know as the "equity premium puzzle" (EPP), see the page of the Federal Reserve Economic Data website myf.red/g/6LsS for a visual representation of the premium in the last 25 years (difference between stock and BBB corporate yields).

Economists have tried to single out ways to explain why people invest far less in stocks than what would be implied by their risk aversion, as measured in other (personal or social) situations. In Benartzi and Thaler [2] it is argued that "the combination of a high equity premium, a low risk-free rate, and smooth consumption is difficult to explain with plausible levels of investor risk aversion" and "myopic loss aversion" is introduced as a possible justification. Barberis et al. [1], in a somewhat similar vein, need (dis)utility from fluctuation of financial wealth. DeLong and Magin [4] survey other approaches, including the use of prospect theory, the value of liquidity or the role of taxation to account for the puzzle.

Agent-based modeling is a methodology to build computational models of real-world systems where autonomous agents (individuals, traders, households, firms, software agents, robots...) interact in various forms, learn, sense the environment and often use fast and frugal heuristics that do not need unrealistic degrees of rationality or processing capability. A good introduction is in Railsback and Grimm [8], that also includes a thorough treatment of the NetLogo programming platform that was used to develop the model presented in this paper. See Steinbacher et al. [11] for a recent review.

Among the features that can be used in agent-based models to analyze possible paths to generate a sizeable equity premium, we investigate the role of direct and indirect interaction among traders. Quite naturally in a financial setup, agents collectively contribute to form the market price of the stock and indirectly affect—and are affected by—the strategies used by the other traders. Moreover, agents are occasionally paired with random peers and contrast the profitability of their strategy, switching to the best-performing one in the quest for improvement. This direct learning scheme is based on pure imitation of successful examples with no need to gather, or elaborate, data or (try to) compute sophisticated conditional equilibria.

The conceptual framework of this paper is inspired by evolutionary game theory. Originally introduced by biologists to analyze with formal tools long-term adaptation of biological populations, the idea that the reproductive fitness depends on the genotypes was later extended to economics, see Sandholm [10] and Newton [7]. Of course, in this setup, agents are not assumed to be genetically pre-programmed but are able to adjust their behaviors or strategies favoring larger payoffs (as opposed to Darwinian fitness). Recently, in Robson and Orr [9] it is argued by means of an evolutionary model that the EPP is due to agents' greater aversion to aggregate risk, such as the one faced in financial markets, with respect to idiosyncratic risk (of more personal nature).

Our agent-based evolutionary model converges to an equilibrium with a over-whelming presence of demand functions (or, if you wish, strategies) which systematically overestimate the variance of the risky asset. As a consequence, a large share of market population hold relatively small amount of stocks. The option to buy a costly information signal to predict return can reduce the effect. However, this holds only if information is accurate, cheap and used in non-volatile markets. Generally speaking, the overestimation of the variance in the long-run, and its related implications, are observed in many of the instances examined in a detailed robustness test. Moreover, in such instances, most of the (survived) equilibrium strategies appear to use little or no information. These results shed a novel light on the puzzle and point to a potential new channel to explain this long-debated anomaly.

Section 2 presents the ABM model, describing the market, the agents and the learning protocol that is naturally used as an evolutionary device to favor trading strategies with higher payoffs. Section 3 presents simulations' results obtained in a benchmark case. Some key parameters are then varied in the following section that shows that results are remarkably robust. We finally conclude with some discussion.

## 2   The Model

This section describes a simple market with a risky stock and a riskless asset. The setup is minimal to keep the focus on the co-evolution of a population of traders who compete for high profits and decide which information and risk factors to take into account in their decisions.

## 2.1   The Market

We assume $N$ heterogeneous agents are given an initial endowment $w_0$ at the beginning of every period, place orders and collect revenues that are immediately consumed at the end of the period. Some agents are then allowed, with some probability, to change their trading strategy using an imitation mechanism that favors the ones with larger revenues. A new population, with a different distribution of strategies, is formed and the game is repeated $T$ times, $t = 1, \ldots, T$.

The riskless asset has unit cost and pays $R = 1 + r$ after one period. There is also a stock in zero net supply with random payoff $\tilde{D}_t = d + \tilde{\theta}_t + \tilde{\epsilon}_t$, where $d$ is a known deterministic component of revenues, $\tilde{\theta}_t$ is an informative signal that can be acquired at a cost of $c$ per period and $\tilde{\epsilon}_t$ is an unobservable noise term (unknown to everyone). We omit $t$ and occurrences of tilde, unless needed, in what follows and notice that $\theta$ can be referred as *information*, as it truly affects the random revenue $D$. Some agents, however, may believe that an uninformative signal $\tilde{\gamma}_t \equiv \gamma$ also affects the outcome. $\gamma$ can be obtained at no cost, if desired, and it can be considered as pure *misinformation* having nothing to do with $D$ (even though agents regard it as helpful).

We assume that $\theta, \epsilon, \gamma$ are normally and independently distributed:

$$\theta \sim N(0, v_\theta), \epsilon \sim N(0, v_\epsilon), \gamma \sim N(0, v_\gamma), \theta \perp \epsilon \perp \gamma,$$

where $v_\theta$, $v_\epsilon$ and $v_\gamma$ are the variances of $\theta$, $\epsilon$ and $\gamma$.

The equilibrium price at any time $t$ is determined by (net) demands of agents solving the equation

$$\sum_{i=1}^{N} x_i(p_t|\mathbf{b}_i, \theta, \gamma) = 0, \tag{1}$$

where $x_i(p_t|\mathbf{b}_i)$ is the demand of the $i$-th agent at price $p_t$ and $\mathbf{b}_i$ is a vector of heterogenous parameters differentiating individual strategic behavior and will be described in detail in the next subsection.

### 2.2  The Agents

The demand of the risky asset is consistent with the idea that agents, as a whole, are aware that $D$ depends on some of (but not necessarily all) the variables mentioned above: for a price $p$ the demand function of the $i$-th agent is

$$x_i(p, \mathbf{b}_i) = \frac{d + b_1\theta + b_2\gamma - pR}{a(b_3 v_\theta + b_4 v_\gamma + v_\epsilon)}, \tag{2}$$

where $\mathbf{b}_i = (b_1, b_2, b_3, b_4)$, $i = 1, \ldots, N$ is a vector of individual bits (i.e., $\mathbf{b}_i \in \{0, 1\}^4$) that can evolve in time due to imitation (and, hence, should be formally denoted as $\mathbf{b}_{it}$ when referring to the $i$-th agent at time $t$). Again, for the sake of exposition, we omit individual and temporal indexes. As $\mathbf{b}$ shapes and determines the trading behavior of the agents, we will refer to it using the term "strategy".

The demand in (2) increases with the perceived average revenue in excess of what would be gained with a riskless investment (see the numerator) and is corrected for the perceived variance, up to the relative risk aversion coefficient $a$, held constant across the population of traders. Each bit $b_j$, $j = 1, \ldots, 4$ can be thought as a way to switch on and off some random variable in Eq. (2). Take, for instance, an agent with strategy $\mathbf{b} = (1, 0, 0, 0)$: she acquires and uses information $\theta$ in the numerator and perceives a residual variance (in the denominator) depending on $\epsilon$ alone. Such an agent can be regarded as *informed*, as she employs $\theta$ and discard $\gamma$, as well as *rational*, as she correctly realizes that the variance of $D$ is not affected by misinformation $\gamma$ or by $\theta$, that is known, but only depends on the noisy and unobservable component $\epsilon$. Indeed, in this paper, rationality refers to the correct understanding of the data-generating process of $D$, and would be achieved in the model when $b_1$ is either 0 or 1, $b_3 = 1 - b_1$ and $b_2 = b_4 = 0$. In other words, a rational agent would ignore

**Table 1** Description of several strategies encoded in the vector $\mathbf{b} = (b_1, b_2, b_3, b_4)$

| Nickname | $b_1$ | $b_2$ | $b_3$ | $b_4$ | Description |
|---|---|---|---|---|---|
| Informed | 1 | 0 | 0 | 0 | *Informed, rational* |
| Prudent | 0 | 0 | 1 | 1 | *Uninformed, irrational, small demand* |
| Uninformed | 0 | 0 | 1 | 0 | *Uninformed, rational* |
| Fearless | 0 | 0 | 0 | 0 | *Uninformed, irrational, relatively large and stable demand* |
| Confused | 1 | 1 | 1 | 1 | *Informed and misinformed, irrational* |

The first two bits are related to the use of the information and misinformation in the prediction of the mean returns; the third and fourth bits are used to compute the perceived risk

bits related to $\gamma$ and would either buy the information or, if not, include it in the denominator of Eq. (2).

By contrast, someone using the strategy $\mathbf{b} = (0, 0, 1, 1)$ may be considered quite *prudent*: indeed, none of the signals $\theta$ or $\gamma$ is used and the perceived variance is large as it includes both the summands $v_\theta$, $v_\gamma$, as well as the ubiquitous $v_\epsilon$. As a consequence, such an agent would trade a much smaller $x$, for any given $p$, than an informed trader. In this specific case, clearly, the strategy is not fully rational as $\gamma$ is erroneously affecting the demand.

Table 1 lists some relevant strategies that can, to some extent, be interpreted in terms of their ability to correctly identify the conditional expected revenue and variance.

While it's not always possible to provide a behavioral description of every strategy encoded in $\mathbf{b}$, Table 1 features a few meaningful examples. For instance, uninformed agents discard useful information (avoiding the cost), but are rational in that they correctly understand the way returns are generated and take into account the uncertainty arising from the unknown $\theta$; traders with a null $\mathbf{b}$, in the last row of the table, are dubbed *fearless* to stress the lack of any risk adjustment in the denominator of (2), an action leading often to relatively large orders.

## 2.3 Learning

At the end of any period $t$, the equilibrium price $p_t$ is computed using Eq. (1). Clearly, $p_t$ is a function of the distribution of the strategies in the population and of realized random variables $\theta_t$, $\gamma_t$, that are known to the agents whose first and second bits are set to 1, respectively. After the unobservable shock $\epsilon_t$ is drawn and uncertainty is resolved, the realized profit for an agent is

$$w_{it} = x_{it} D_t + (w_0 - x_{it} p_t) R - c \cdot b_1,$$

where the first component is the revenues arising from $x$ units of the risky stock, the second part comes from investing in the riskless asset all the cash that was not used to get the stocks, and $cb_1$ is the cost of getting the information.

Learning is based on agents' pairwise comparisons of the profits. In detail, we form $h < N/2$ random couples of traders and, letting individuals $i$ and $j$ be one such couple, the vectors **b** are updated using:

$$\text{If } w_{it} > w_{jt}, \ \mathbf{b}_{j,t+1} = \mathbf{b}_{i,t};$$
$$\text{If } w_{it} < w_{jt}, \ \mathbf{b}_{i,t+1} = \mathbf{b}_{j,t};$$
$$\text{If } w_{it} = w_{jt}, \ \text{no change}.$$

The interpretation of this learning scheme is immediate in terms of evolutionary game theory: agents using possibly different strategies obtain different payoffs; they occasionally meet another peer and revise their strategy switching to the one with bigger profits (pure imitation of better strategies); as a consequence, strategies with better payoffs tend to increase their relative frequency (which, in turn, may alter their future success).

Slightly more formally, a population $\mathcal{B}_t = \{\mathbf{b}_{it} : i = 1, \ldots, N\}$ of agents (or strategies) at time $t$ evolves using the above revision protocol to obtain $\mathcal{B}_{t+1}$, which has at most $h$ differences with respect to $\mathcal{B}_t$. The relative frequencies of each of the 16 strategies are then investigated letting $t$ reach $T$, for large $T$.

## 3 Results

This section discusses the results obtained simulating in NetLogo [12] the agent-based model described in the previous section.[1] We first illustrate the outcomes in a benchmark case and then show how results change varying systematically some of the parameters of the model.

### 3.1 The Benchmark Case

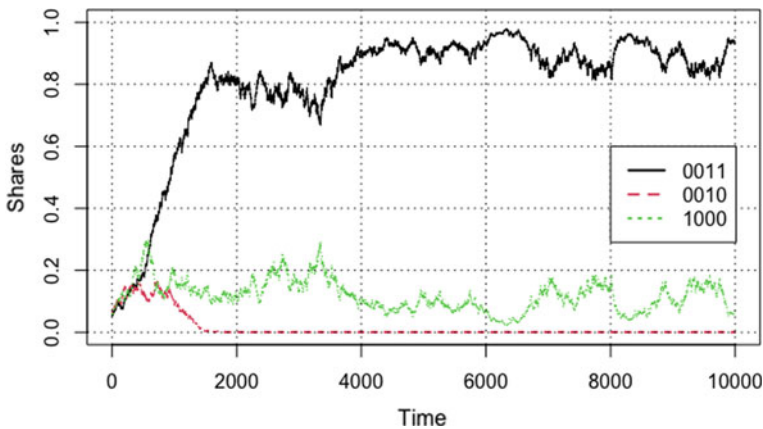Table 2 lists the values taken by the parameters of the model in a benchmark configuration.

The values are roughly representative of a market where, for instance, one period is one year, the riskless rate is 1%, the standard deviation of the revenues of the risky asset is 20% = $\sqrt{0.04}$, the standard deviation of information (and misinformation) is 10% and the cost of acquiring the informative signal is 3% (that may be a reasonable approximation of the fees of a financial professional providing valuable advice).

---

[1] The code is available on the website of the authors.

**Table 2** Values of the parameters of the model used in the benchmark case, with a brief description

| Param. | Value | Description | Param. | Value | Description |
|---|---|---|---|---|---|
| $N$ | 1000 | Number of agents | $R$ | 1.01 | Gross return of riskless asset |
| $T$ | 10,000 | Trading periods | $a$ | 2 | Risk aversion coefficient |
| $c$ | 0.03 | Cost of information $\theta$ | $v_\theta$ | 0.01 | Variance of information $\theta$ |
| $h$ | 15 | Learning couples | $v_\epsilon$ | 0.04 | Variance of noise $\epsilon$ |
| $d$ | 1.1 | Part of stock revenue | $v_\gamma$ | 0.01 | Var. of misinformation $\gamma$ |



**Fig. 1** Time series of the shares of strategies 0011 (prudent), 1000 (informed) and 0010 (uninformed) in a standard simulation run of the benchmark case lasting 10,000 periods. In particular, $c = 0.03$, $v_\epsilon = 0.04$

Figure 1 shows how the fractions of informed (1000), prudent (0011) and uninformed (0010) agents evolve in 10,000 periods in one standard simulation run.[2] It can be seen that about 2000 periods suffice to reach a homeostatic equilibrium where the share of prudent traders hovers around 90%, informed agents are 10% and we observe the extinction of the uninformed (as well as any **b** other then with 1000 and 0011.

The result that Strategies 1000 and 0011 are the only survivors in the long run is a first and important regularity of the model for this parameters' choice. Figure 2 is based on 100 simulations (of 10,000 periods) and depicts the mean share of all the strategies (the box) equipped with standard deviations (equal to the length of the vertical line extending over the bars). On average, equilibrium is reached when about nine tenths of agents are prudent and the remaining ones are informed.

Quite remarkably, Fig. 2 shows that when equilibrium is reached evolutionary pressure has obliterated all other strategic variations. Notice that, at equilibrium, the

---

[2] All simulations are initialized setting the bits in **b** randomly in {0, 1}.

**Fig. 2** Average shares and related standard deviation of different strategies in the benchmark case (with $c = 0.03$, $v_\epsilon = 0.04$). Mean values are shown based on 100 simulations of 10,000 periods. Essentially, only strategies 0011 and 1000 survive

probability that the profits of the informed are bigger than the ones of the prudent is 50% (this holds because, essentially, learning forces the surviving strategies to have the same median profits and, if this were not the case, the shares would have drifted away from that equilibrium in the presence of a tendency to prefer one of the two strategies).

This outcome suggests a novel explanation of the EPP from the bottom up. The large majority of prudent traders underinvest in the risky asset, being their demand particularly low as observed in the previous section. Indeed, this is due to a systematic overestimation of the variance of the gains of the stock that, in turn, reduces the demand of the risky asset and favors more conservative savings in the safe bond. Even if the prudent strategy reduces the average profit of investment, nonetheless the median revenues of the prudent are the same as the ones of the informed agents (at the end of each period, when consumption takes place).

Put differently, traders demanding small amounts of the risky asset become very popular in a market where they "compete" according to the (sharp) rule described in Sect. 3.2 and occasionally compare their profits, achieving a reduction of the risk of being pushed out of the market in the long run. Such a majority of prudent traders fits very well the puzzling observation that fewer agents than expected invest in equities. Assuming an unrealistic level of sophistication, it could perhaps be argued that agents would maximize utility or realize that larger mean profits can be traded for the smaller gains obtained in about 50% of time. However, simple and straightforward comparisons based on the question "does $\mathbf{b}_i$ produce larger gains than $\mathbf{b}_j$?" are strong calls to immediate action and more convincing behavioral drivers.

**Table 3** Modal strategy at equilibrium, for several values of parameters $c$ and $v_\epsilon$

| $v_\epsilon$ | Cost $c$ | | | | |
|---|---|---|---|---|---|
|  | 0.01 | 0.02 | **0.03** | *0.04* | 0.05 |
| 0.03 | 0.54/1000 | 0.66/0011 | 0.80/0011 | 0.92/0011 | 0.90/0011 |
| **0.04** | 0.50/0011 | 0.73/0011 | **0.92/0011** | 0.75/0011 | 0.34/0000 |
| 0.05 | 0.59/0011 | 0.84/0011 | 0.78/0011 | 0.37/0011 | 0.51/0000 |
| *0.06* | 0.63/0011 | 0.93/0011 | 0.50/0011 | *0.51/0000* | 0.61/0000 |

The entry $s/\mathbf{b}$ means that the modal share $s$ was reached by agents using strategy $\mathbf{b}$. The boldfaced entry is relative to the benchmark configuration and the italicized one is discussed in the text

## 3.2 Robustness Tests

As expected, the outcomes described previously are sensitive to the parameters of the market. In this subsection, we explore the robustness of the results with respect to changes in the cost $c$ of information and in the size $v_\epsilon$ of the unobservable and idiosyncratic shock. We use BehaviorSpace, a NetLogo's tool that allows to run experiments and gather data systematically "sweeping" (portions of) the parameters' space. Table 3 shows, for $c \in \{0.01, \ldots, 0.05\}$ and $v_\epsilon \in \{0.03, \ldots, 0.06\}$, the largest average share at equilibrium, based on 100 simulations of 10,000 periods for each of the 20 couples $(c, v_\epsilon)$. For instance, corresponding to the benchmark parameters (boldfaced in Table 3), we see that the largest share (92%) is made of prudent investors (with $\mathbf{b} = 0011$).

Table 3 demonstrates that prudent investors dominate the scene in many cases, with shares varying from 37 to 93%. Therefore, to a great extent, quite some under-investment in stocks is natural. Informed agents are prevalent at equilibrium only when the cost of information and the variance of the noise are low, in the top-left corner of the Table where $c = 0.01$ and $v_\epsilon = 0.03$.

Interestingly, the bottom-right corner of Table 3 shows that the majority share at equilibrium is made of *fearless* agents, as they are nicknamed in Table 1. When both the cost of information and the variance of $\epsilon$ are large, a fraction of agents thrive with no use of signals $\theta$ or $\gamma$ and avoiding any adjustment for the variance of the risky profits. Figure 3 depicts the equilibrium shares in one of such market instances, when $c = 0.04$, $v_\epsilon = 0.06$.

Only strategies 0000, 0001, 0010, 0011 survive in the long run in this market where information is more costly and market returns are more volatile. This is a plausible explanation of why Strategy 1000 died out but, truly, a reason for the observation that no survivor switches on bit $b_1$, that would imply information is bought and used, or $b_2$, that would imply misinformation is used (indeed, in this case we have that $b_1 = b_2 = 0$ for all agents, excluding a handful of outlying traders using 0110 or 0111, see Fig. 3). Such a market is populated by many individuals who do not use any information (or misinformation, for what it matters), and who keep at 1 at least one of the bits $b_3$, $b_4$ located in the denominator of Eq. (2), reducing on average the quantity of stock kept in their equilibrium portfolio.
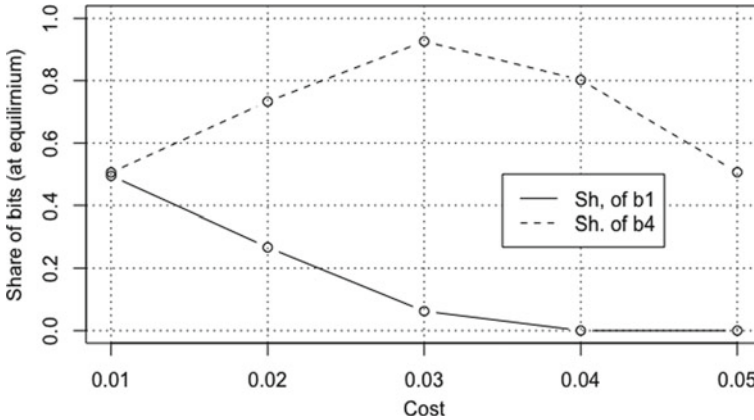
**Fig. 3** Average shares and related standard deviation of different strategies when $c = 0.04$, $v_\epsilon = 0.06$. Mean values are based on 100 simulations of 10,000 periods. Only Strategies 0000, 0001, 0010 and 0011 survive (with tiny exceptions), fearless traders are over 50% and very few agents use information, as discussed in the text

The main conclusion of our robustness test is that agents invest far less in the stock market than would be implied by a full-fledged (and probably unrealistic) model of rational allocation. This holds even under variations of several key parameters of the model (that, in some cases, lead altogether to the disappearance of information from the strategies).

## 4 Discussion and Conclusion

One of the most interesting features of ABMs is their ability to accommodate for heterogenous features of the agents. We have considered a bunch of strategies that differ in the content of information that is used to assess expected profits (in the numerator of Eq. 2), as well as in the risk factors that are considered (in the denominator). Even though, in principle, agents could use any combination of active bits, evolutionary pressure wipes out most of the strategies. As pointed out by a reviewer, whether only one strategy survives in the very long run is still an open question (and 10,000 periods may not be enough to reveal the steady state). Further research should also consider the effects of the introduction of mutation or the replacement of some agents with new ones. In any case, many strategies may be assumed to be of little importance asymptotically. In particular, in the benchmark case only informed and prudent traders survive; with the latter keeping a prominent position in many other situations. Prudent and fearless strategies do not use information $\theta$ and this results in a small number of information users, as shown in the solid line of Fig. 4 depicting

**Fig. 4** Shares of agents with $b_1 = 1$ (solid line) and $b_4 = 1$ (dashed line), as a function of the cost $c$ when $v_\epsilon = 0.04$. The former make some use of the information and the latter take misinformation $\gamma$ as a risk factor

the share of traders whose $b_1 = 1$ at equilibrium for $v_\epsilon = 0.04$. A similar tendency is reported in Gerotto et al. [5] in a setup with only two strategies.

The frequency of $b_4 = 1$ at equilibrium is even more relevant to explain the EPP and can be associated to a majority of traders who buy limited amounts of stocks because their strategy inflates the perception of risk and reduce traded quantities. The dashed line in Fig. 4 shows the share of users with $b_4 = 1$: with the exception of a few markets where fearless traders prevail, most agents are extremely cautious for all levels of cost when $v_\epsilon = 0.04$ (and this fraction is substantial and rarely falls below 50% in the many parametric combinations we have investigated in Table 3).

Overall, our model suggests that the EPP stems, to some extent, from the evolutionary updating of strategies used by myopic traders, where the myopia mainly lies in the assumption that learning is performed with an eye on one-period performance only. While the introduction of long-term orientation may reduce the effect, the salience of recent rewards is well-documented, see Cosemans and Frehen [3] for a recent treatment, and can, in combination with strategy-switching, help in clarifying some of the issues raised by the EPP.

# References

1. Barberis, N., Huang, M., Santos, T.: Prospect theory and asset prices. Q. J. Econ. **116**(1), 1–53 (2001). http://www.jstor.org/stable/2696442
2. Benartzi, S., Thaler, R.: Myopic loss aversion and the equity premium puzzle. Q. J. Econ. **110**(1), 73–92 (1995)
3. Cosemans, M., Frehen, R.: Salience theory and stock prices: empirical evidence. J. Finan. Econ. **140**(2), 460–483 (2021)
4. DeLong, J.B., Magin, K.: The U.S. equity return premium: past, present, and future. J. Econ. Perspect. **23**(1), 193–208 (2009). http://www.jstor.org/stable/27648300
5. Gerotto, L., Pellizzari, P., Tolotti, M.: Asymmetric information and learning by imitation in agent-based financial markets. In: De La Prieta, F., González-Briones, A., Pawleski, P., Calvaresi, D., Del Val, E., Lopes, F., Julian, V., Osaba, E., Sánchez-Iborra, R. (eds.) Highlights of Practical Applications of Survivable Agents and Multi-Agent Systems. The PAAMS Collection, pp. 164–175. Springer International Publishing, Cham (2019)
6. Mehra, R., Prescott, E.C.: The equity premium: a puzzle. J. Monetary Econ. **15**(2), 145–161 (1985)
7. Newton, J.: Evolutionary game theory: a renaissance. Games **9**(2) (2018). https://doi.org/10.3390/g9020031, https://www.mdpi.com/2073-4336/9/2/31
8. Railsback, S., Grimm, V.: Agent-Based and Individual-Based Modeling: A Practical Introduction, 2nd ed. Princeton University Press (2019). https://books.google.it/books?id=X3SYDwAAQBAJ
9. Robson, A.J., Orr, H.A.: Evolved attitudes to risk and the demand for equity. Proc. Natl. Acad. Sci. **118**(26), e2015569118 (2021)
10. Sandholm, W.H.: Population Games and Evolutionary Dynamics. Economic Learning and Social Evolution. MIT Press (2010)
11. Steinbacher, M., Raddant, M., Karimi, F., Camacho Cuena, E., Alfarano, S., Iori, G., Lux, T.: Advances in the agent-based modeling of economic and social behavior. SN Bus. Econ. **1**(7), 99 (2021)
12. Wilensky, U.: NetLogo. Center for Connected Learning and Computer-Based Modeling, Northwestern University, Evanston, IL (1999). http://ccl.northwestern.edu/netlogo/

# Energy and Climate

# An Agent-Based Model of UK Farmers' Decision-Making on Adoption of Agri-environment Schemes

Chunhui Li, Meike Will, Nastasija Grujić, Jiaqi Ge, Birgit Müller, Arjan Gosal, and Guy Ziv

**Abstract** Agri-environment schemes (AES) are government-funded voluntary programs that incentivise farmers and land managers for environmental friendly farming practices. Understanding farmers' decision-making process and its impact on AES adoption can aid policy makers in designing AES schemes that meet adoption goals and environmental targets. Farmers' decision-making is complex and involves a range of social, behavioural, economic and ecological factors. In this paper, we present a spatially explicit agent-based model (ABM)—BESTMAP-ABM-UK that simulates farmers' decision-making process, inclusive of farmers' social, behavioural and economic factors, on adopting buffer strips, cover crops, grassland management and arable land conversion to grassland schemes in the UK. The model produces farmers' AES adoption under varied AES scheme designs in term of the contract length, the offered payment level, the bureaucracy level and the required minimal area. We apply the Morris screening method to analyse the importance of the model parameters in a status quo scenario, in which current UK AES designs are used. The results show that the average accepted payments of farmers for buffer strips and grassland management and farmers' intrinsic openness to buffer strips have the most significant impact on the farm adoption rate in the model.

**Keywords** Agri-environment schemes · Agricultural agent-based model · Decision-making

C. Li (✉) · J. Ge · A. Gosal · G. Ziv
University of Leeds, Leeds LS2 9JT, UK
e-mail: jenny.c.li@outlook.com; c.li2@leeds.ac.uk

M. Will · B. Müller
Helmholtz Centre for Environmental Research—UFZ, 04318 Leipzig, Germany

N. Grujić
BioSense Institute, University of Novi Sad, 21000 Novi Sad, Serbia

# 1　Introduction

Agri-environment schemes (AES) are government-funded voluntary programs to incentivise farmers and land managers for environmental friendly farming practices. Understanding farmers' decision-making process in AES adoption can aid policy makers in designing AES schemes that meet adoption goals and environmental targets.

Farmers' decisions-making on AES adoptions is complex and involves a range of social, behavioural, economic and ecological factors [1, 2]. A farm's biophysical and structural characteristics affect the farmer's options and decisions. Large sizes and wealthier farms participate more in AES because they are more likely to have suitable lands for conservation practices and sufficient resources of labour and finance means [2–9]. Farm types, location, and soil quality of farm lands impact on the AES adoption [2, 4, 5, 10–13] as well as demographics of farmers, i.e., farmers' age, education and work arrangement (i.e., full-time or part-time farmers) [2, 3, 5–7, 11]. Farmers' social capital, e.g., connections with advisory services, family and friends networks, and local communities, is an influential factor [2, 6, 10]. In addition, social and behavioural factors, e.g., farmers' previous experience and knowledge about AES [3, 4, 11], personal attitude towards environment [2, 6] and schemes' effectiveness [14], and trust in agriculture authorities [2, 15] were also found to be important.

ABM is a useful tool to study complex systems in various research fields. A significant number of ABMs have been developed to model different aspects of the agricultural system. Extensive literature of ABMs can be found covering farm management, land use, agricultural economy, agricultural policy evaluation, environmental change, climate change, and a mixture of these subjects [16–21]. Existing models that study impact of agricultural policies emphasise economic aspects of farmers' decisions, for example, AgriPoliS [18] and IFM-CAP [22], while farmers' emotions, values, learnings and social interactions receive little attention in the models [23].

In the following sections, we describe the spatial-explicit BESTMAP-ABM-UK that models UK farmers' decision-making on AES adoptions inclusive of farmers' social, behavioural and economic factors in a three-step decision framework. BESTMAP-ABM-UK is one of the BESTMAP-ABMS that are designed to study farmers' decision-making on AES adoptions in five European regions. BESTMAP-ABMS is a part of the BESTMAP policy impact assessment model (PIAM) suite to link economic modelling, behavioural ABM and established ecosystem service modelling [24]. In BESTMAP-ABM-UK we use the Humber region in England as the case study area. The Humber region has approximately 3500 farms and 56,000 agricultural parcels covering about 350,000 hectares (ha) in the area, according to Rural Payments Agency 2019 data. We carry out global sensitivity analysis using the Morris screening method. The Morris screening method is a computationally efficient screening technique that allows us to identify the important model input factors. The results show the importance rank of seventeen input parameters, including farmers' mean accepted payments and the standard deviation, openness

**Table 1** AES typologies in the ABM

| AES types | CSS option codes |
|---|---|
| Buffer strips (BS) | SW1, SW2, SW3, SW4, SW11, AB1, AB3, AB8, WT2 |
| Grassland management (GM) | GS2, GS5, GS6, GS7, GS9 |
| Cover crops (CC) | SW6 |
| Arable land conversion to grassland (CVN) | SW7 |

related factors including prior experience, intrinsic openness, influence from advisory and social network, and farmers' preferred amount of AES area in their farms.

## 2 Model Design

### 2.1 AES Scheme Options in the ABM

A wide variety of AES is designed targeting specific farming systems and ecosystems in different states across the EU. Coordinating with five case studies within the BESTMAP project, we select four types of AES: buffer strips, grassland management, cover crops, and arable land conversion to grassland, according to the relative importance in terms of spatial coverage and findings in the interviews with the farmers in the five case study regions [25]. In the UK case study, we select relevant scheme options according to the scheme characteristics and group them into the four types out of hundreds of scheme options in the countryside stewardship scheme (CSS). Table 1 lists the CSS option codes that belong to each AES group. More details of the AES options can be found on the UK government website.[1]

### 2.2 Agents and Agent Behaviours

**Agents and their state**. Farmers, fields and AES contracts are modelled as agents.

- Farmer-agents represent land managers managing a set of fields and making decisions on farm operations. A farmer-agent has these attributes: ID, location, total agricultural area, farm fields, farm system archetype (FSA), economic size, openness towards a specific AES that is derived from farmers' previous AES experience, intrinsic openness, influences from advisory service and social networks that it belongs to, the willingness to accept (WTA) and the envisioned area for

---

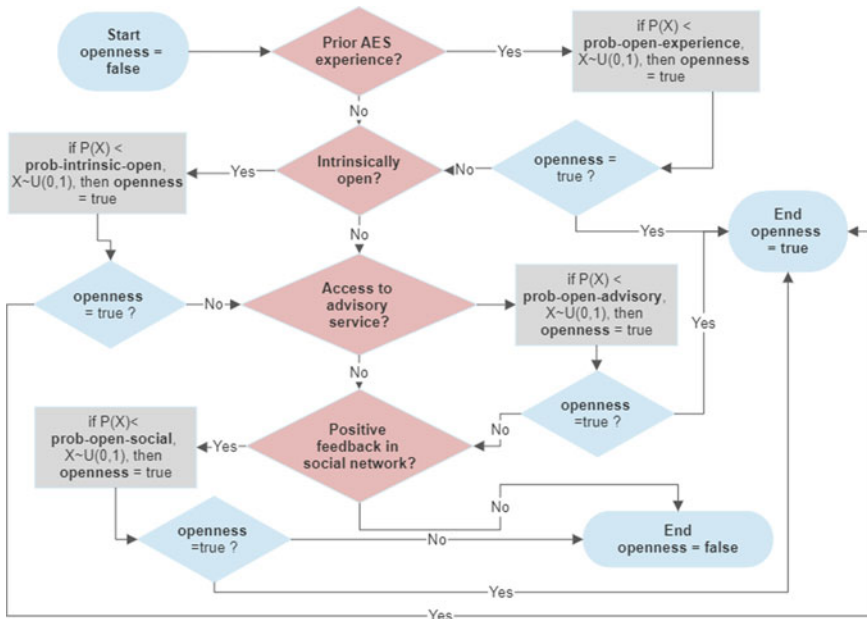[1] https://www.gov.uk/countryside-stewardship-grants.

each type of AES. FSAs in BESTMAP project are characterised into five types: general cropping, horticulture, permanent crops, grazing livestock and mixed farms, assuming farmers operating similar farms respond to policies in a similar fashion. The economic sizes of farms are large, medium, small, and smaller than 2000 euros. WTA refers to the minimum monetary value that a person is willing to accept for selling a service or bearing some harm in economics [26]. The WTAs in our model will be derived from the discrete choice experiment (DCE) we are carrying out. A farmer-agent's WTA is influenced by the required contract length, the bureaucracy level and the minimum area of AES.

- Fields-agents represent agriculture parcels that are managed by farmer-agents. They are spatial polygons that have attributes including parcel ID, location, field size, land use type, owner and soil quality. The land use types include arable, grassland, horticulture, woodland and others.
- AES-agents record the AES contract information including contract ID, farmer participant, parcel ID and the contract area. AES agents are endogenous in the model and generated at runtime following farmers' decision on taking up AES.

**Agent networks**. Farmers who have the same FSA are connected in the region.

**Farmer-agent behaviours**. Farmer-agents' decision-making takes three steps.

- In the first step, farmer-agents check whether they are open to a particular type of AES. The openness is derived from the previous experience, the intrinsic openness, and the influence from advisory services and social networks. The flow chart in Fig. 1 depicts this process. If farmer-agents are not open to any of the four AES, they exit the decision-making process; if they are open to any of the AES, they proceed to the next step.
- In the second step, farmer-agents check whether they have suitable lands for the AES that they are open to adopt. The selection is based on land use types: arable lands are eligible for buffer strips, cover crops and conversion of arable land to grassland AES; Grasslands are eligible for buffer strips and grassland management.
- In the third step, farmer-agents deliberate on which AES to adopt and the details of the location and the contract area. If farmer-agents are open to one or more AES options and have suitable fields for them, they compare the offered payment (in pounds per hectare) with their WTA and choose an AES based on the profit. When a farmer-agent decides to take up an AES, the contract details including the location and area is decided and an AES-agent is generated. Farmer-agents tend to put smaller and less productive fields, i.e., fields with lower soil organic matter, under AES. The AES contract area is set based on the field-level adoption proportion pattern in historic CSS adoption data. On average, buffer strips AES takes 6.8% of a field, cover crops AES takes 63.6% of a field, grassland management takes 96.1% of a field and arable land conversion to grassland AES takes 82.5% of a field.

**Fig. 1** The first step of decision-making process—deriving openness. Farmer-agents' openness status (true or false) is decided by four factors—whether a farmer-agent has prior experience, whether it is intrinsically open to the four types of AES, whether it has access to advisory service, whether other farmer-agents in its social network has positive experience. The impact of the four factors are set to be four probabilities: "prob-open-experience", "prob-intrinsic-open", "prob-open-advisory" and "prob-open-social", subject to standard uniform distribution U(0,1). At the end of this process, the farmer-agent either goes into the next step (if openness is true) or exits the decision-making process (if openness is false)

## 2.3 Model Data Sources and Stochasticity

**Data sources**. The model is parameterised using four sets of data. Firstly, we use Land Parcel Identification System (LPIS) data received from UK Rural Payments Agency to set up farmer- and field- agents. Secondly, CSS adoption data from Natural England[2] is used to inform the model baseline of AES adoption in term of adoption volumes and patterns. Thirdly, we use the mean estimates of carbon density in topsoil published by UK Centre for Ecology & Hydrology[3] to set soil quality parameter of field agents. The last, we conduct a DCE in Humber aiming to study farmers' preferences and the WTAs under different scheme design conditions.

---

[2] https://naturalengland-defra.opendata.arcgis.com/datasets/Defra::countryside-stewardship-scheme-agreements-england/about.

[3] https://eip.ceh.ac.uk/naturalengland-ncmaps/reportsData.

**Model stochasticity**. In deriving famer-agents' openness, we apply a standard uniform distribution to set the influence of the modelled factors. Therefore, a farmer-agent's openness status can vary at different simulation time ticks. While waiting for the DCE to complete, we set random values for farmer-agents' WTA assuming that the WTA towards one AES in farmer population is subject to the normal distribution with a mean value and a given standard deviation. The WTA of farmer-agents is set in the initialisation and stays fixed in a simulation. Field-agents in the model are static agents.

## 2.4 *The Model Implementation*

The model is implemented in NetLogo [27]. Simulations and sensitivity analysis are run using the R package nlrx [28].

## 3 Preliminary Results and Discussion

### 3.1 *Global Sensitivity Analysis*

We run the Morris screening method, a.k.a. Morris elementary effect screening method, for global sensitivity analysis, with the aim of better understanding the importance of the parameters in the model. Elementary effects refer to the changes of model output that are solely due to changes in a particular model input [29]. Morris screening is a one-factor-at-a-time (OAT) method that produces two sensitivity measures based on the elementary effects of model input factors' samples: one is the estimated mean of elementary effects, measuring the overall effect of a parameter on the output, noted as $\mu$; the other is estimated standard deviation indicating a dependency on other parameters, noted as $\sigma$. A revised measure $\mu^*$ is "the estimate of the mean of the distribution of the absolute values of the elementary effects" [30]. Although $\mu^*$ is sufficient to produce a ranking of the parameters, this measure does not show the positive or negative effect of a parameter on the model. Therefore, we use all three measures in our analysis. Detailed algorithm description of the Morris screening method can be found in the book by Saltelli et al. [31].

### 3.2 *Experiment Design*

We set our model to run the status quo scenario, in which we set the AES designs to represent the current CSS designs in England. We run the model for one tick

**Table 2** AES parameter values in the ABM

| AES types | Payment (£/ha) | Length (years) | Bureaucracy | Minimal area (ha) |
|---|---|---|---|---|
| Buffer strips (BS) | 524 | 5 | Medium | 0 |
| Grassland management (GS) | 124 | 5 | Medium | 0 |
| Cover crops (CC) | 183 | 5 | Medium | 0 |
| Arable land conversion to grassland (CVN) | 321 | 5 | Medium | 0 |

(representing one year), in which farmers' prior AES experiences are set according to the CSS AES adoption data in 2019 at the beginning of a simulation.

The AES scheme parameter values are listed in Table 2. The payment is the average payment level of the AES in each AES group in Table 1. As the minimal area is not explicitly stated in the CSS descriptions, it is set to be 0 in the experiment. We assume current schemes are medium bureaucracy level for farmers.

**The input data sample**. We run the experiment and analysis based on a sample data containing 352 farmers to reduce the computational cost in the experiment. The farmers are sampled based on the farmers' FSA and account for 10% of the whole Humber region farmers. We compare the proportions of different types of fields in the whole dataset and in the sample data (Fig. 2) and conclude that the sample data is representative of the whole dataset, because the proportions of arable fields and grassland fields in the whole dataset and the sample dataset are at similar levels and these two types of fields are the eligible land types for the four modelled AES.
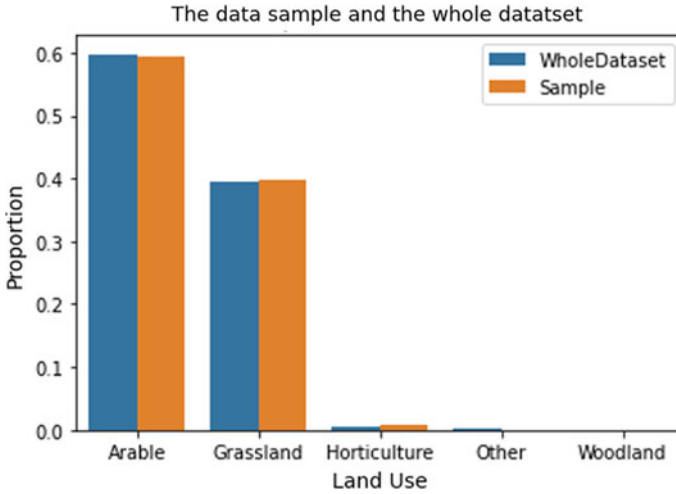
**Model output measure**. The model output is measured by total farm adoption rate. The total farm adoption rate $R_{\text{farm}}$ is calculated by the number of farms participating AES $N_{\text{AES}-\text{farms}}$ divided by total number of farms in the region $N_{\text{farms}}$:

$$R_{\text{farm}} = N_{\text{AES}-\text{farms}}/N_{\text{farms}} \tag{1}$$

## 3.3 Results and Discussion

The Morris screening results are listed in Table 3. The model is run for 50 times for each parameter setting and the mean values of Morris measures $\mu$, $\mu^*$, and $\sigma$ are calculated for discussions in this section.

We tested seventeen model parameters, including: the average of WTA for buffer strips (WTA-BS), cover crops (WTA-CC), grassland management (WTA-GM) and conversion of arable land to grassland (WTA-CVN) in the farmer-agent population, the standard deviation of the WTAs (SD-WTA), the envisioned areas for the four types of AES (envisioned-area-BS, envisioned-area-CC, envisioned-area-GM and

The data sample and the whole datatset

**Fig. 2** The fraction of number of fields in different land use in Humber region. Arable fields, grassland fields and horticultural fields account for 59.8%, 39.6% and 0.6% respectively among all agriculture fields. In the sample data, arable fields and grassland fields account for 59.4, 39.8 and 0.8%

**Table 3** Morris screening results

| Parameters | $\mu^*$ | $\sigma$ | $\mu$ |
|---|---|---|---|
| WTA-BS | 0.237 | 0.204 | −0.237 |
| WTA-GM | 0.17 | 0.279 | −0.17 |
| WTA-CC | 0.072 | 0.119 | −0.072 |
| WTA-CVN | 0.105 | 0.125 | −0.105 |
| SD-WTA | 0.103 | 0.04 | −0.103 |
| prob-intrinsic-open-BS | 0.165 | 0.294 | 0.165 |
| prob-intrinsic-open-CC | 0.556 | 0.072 | 0.556 |
| prob-intrinsic-open-GM | 0.057 | 0.075 | 0.055 |
| prob-intrinsic-open-CVN | 0.110 | 0.124 | 0.115 |
| access-to-advisory | 0.073 | 0.073 | 0.070 |
| prob-open-advisory | 0.073 | 0.129 | 0.073 |
| prob-open-social | 0.070 | 0.094 | 0.070 |
| prob-open-experience | 0.015 | 0.022 | 0.0005 |
| envisioned-area-BS | 0 | 0 | 0 |
| envisioned-area-CC | 0 | 0 | 0 |
| envisioned-area-GM | 0 | 0 | 0 |
| envisioned-area-CVN | 0 | 0 | 0 |

envisioned-area-CVN), the probability that a farmer-agent has access to advisory services (access-to-advisory) and the probabilities that a farmer-agent becomes open because of the advisory services (prob-open-advisory), its prior AES experience (prob-open-experience), the social network (prob-open-social) and being intrinsically open to the four types of AES (prob-intrinsic-open-BS, prob-intrinsic-open-CC, prob-intrinsic-open-GM and prob-intrinsic-open-CVN).
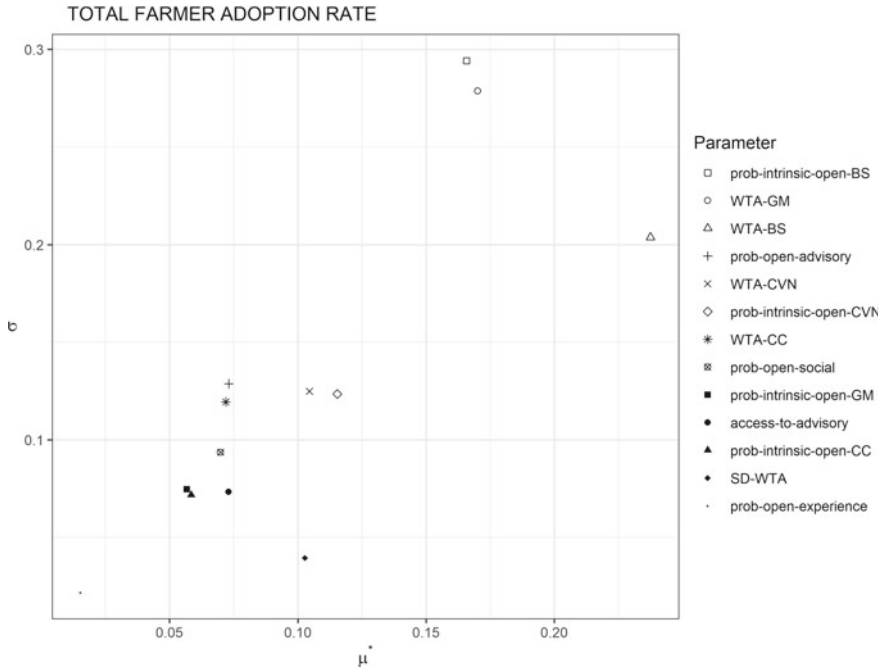
The $\mu^*$ value indicates the relative importance of a parameter to the model output. Therefore, we categorise the parameters' importance into four levels—high, medium, low and minimum according to $\mu^*$ (Table 4). The parameters of high importance include the average WTA for buffer strips, the average WTA for grassland management and the probability of a farmer-agent being intrinsically open to buffer strip, among which the average WTA of buffer strips is the most influential parameter to the total farm adoption rate. The parameters of medium importance include the average WTA for cover crops and arable land conversion to grassland, the probability of a farmer-agent having access to advisory, the probability of a farmer-agent becoming open due to advisory services, and the probabilities of a farmer-agent being intrinsically open to grassland management, cover crops and arable land conversion to grassland. Envisioned areas for the four types of AES have minimum impact on the total farm adoption rate, because envisioned areas do not play a role in a farmer's decision on adopting or not, but affect the amount of contract area once it decides to participate.

We notice that the probability of farmer-agents' being open due to previous experience falls into the low importance category. Because farmer agents' prior experience is updated every tick in a simulation, it has a cumulative effect when running the model for more ticks. Therefore, the importance of the prior experience factor compared to other openness-related factors is underestimated in this one-tick simulation setting.

The negative $\mu$ values indicate that the farmer-agents' average WTA has negative impact on the total farm adoption rate. In comparison, the positive $\mu$ values indicate that the openness-related parameters have a positive impact on the model output measure.

**Table 4** Importance of the model parameters

| Importance | Mean of $\mu^*$ | Parameters |
|---|---|---|
| High | > 0.15 | WTA-BS, WTA-GM, prob-intrinsic-open-BS |
| Medium | > 0.05 and < 0.15 | WTA-CC, WTA-CVN,, SD-WTA, access-to-advisory, prob-open-advisory, prob-open-social, prob-intrinsic-open-CC, prob-intrinsic-open-GM, prob-intrinsic-open-CVN |
| Low | < 0.05 and > 0 | prob-open-experience |
| Minimum | 0 | envisioned-area-BS, envisioned-area-CC, envisioned-area-GM, envisioned-area-CVN |

**Fig. 3** The Morris sensitivity measures μ* and σ for the thirteen influential parameters in the model on total farm adoption rate

We plot the rest parameters in the (μ*, σ) plane suggested by Saltelli et al. [31] in Fig. 3, excluding the minimum important parameters (i.e., envisioned areas). The figure shows that the input parameters with higher μ* values also have higher σ values, which implies that the effect of these parameters on the model output is non-linear.

The Morris method is useful in the model verification and carlibration. On one hand, the results demonstrate the linkage between our conceptual model and the computer implementation, as a verification of the model; On the other hand, the analysis offers us insights about the model parameters' importance on the specified model output, which can be used to prioritise the most influential factors in the carlibration.

## 4 Conclusion and Future Work

In this paper, we introduced the BESTMAP-ABM-UK, which simulates Humber farmers' decision-making process in consideration of farmers' social, behavioural and economic factors when deciding the participation of buffer strips, cover crops, grassland management and arable land conversion to grassland. We applied the

Morris screening method as a global sensitivity analysis in running a status quo scenario. The results show that the farmers' average WTA for buffer strips and grassland management and the probability of a farmer being intrinsically open to buffer strips affect the total farm adoption rate most significantly. The higher the farmers' average WTA is, the less they adopt, while the higher the likelihood of farmers' being open, the more farmers adopt AES, which confirms the conceptual design of the model. These findings not only verify that the computer model is an accurate implementation of the conceptual design, but also inform us the most influential parameters in the model.

In the future, we will apply the Morris screening method in other experiments, including running the model for multiple ticks and using other model output measures, for example, the total number of AES contracts and adoption prediction error. We will then prioritise our focus on the most influential parameters in the model calibration stage.

# References

1. Uthes, S., Matzdorf, B.: Studies on agri-environmental measures: a survey of the literature. Environ. Manage. **51**, 251–266 (2013). https://doi.org/10.1007/s00267-012-9959-6
2. Lastra-Bravo, X.B., Hubbard, C., Garrod, G., Tolón-Becerra, A.: What drives farmers' participation in EU agri-environmental schemes?: results from a qualitative meta-analysis. Environ. Sci. Policy **54**, 1–9 (2015). https://doi.org/10.1016/j.envsci.2015.06.002
3. Pavlis, E.S., Terkenli, T.S., Kristensen, S.B.P., Busck, A.G., Cosor, G.L.: Patterns of agri-environmental scheme participation in Europe: indicative trends from selected case studies. Land Use Policy **57**, 800–812 (2016). https://doi.org/10.1016/j.landusepol.2015.09.024
4. Wilson, G.A., Hart, K.: Financial imperative or conservation concern? EU farmers' motivations for participation in voluntary agri-environmental schemes. Environ. Plan. A **32**, 2161–2185 (2000). https://doi.org/10.1068/a3311
5. Zimmermann, A., Britz, W.: European farms' participation in agri-environmental measures. Land Use Policy **50**, 214–228 (2016). https://doi.org/10.1016/j.landusepol.2015.09.019
6. Siebert, R., Toogood, M., Knierim, A.: Factors affecting European farmers' participation in biodiversity policies. Sociol. Ruralis **46**, 318–340 (2006). https://doi.org/10.1111/j.1467-9523.2006.00420.x
7. Ruto, E., Garrod, G.: Investigating farmers' preferences for the design of agri-environment schemes: a choice experiment approach. J. Environ. Plan. Manag. **52**, 631–647 (2009). https://doi.org/10.1080/09640560902958172
8. Hynes, S., Garvey, E.: Modelling farmers' participation in an agri-environmental scheme using panel data: an application to the rural environment protection scheme in Ireland. J. Agric. Econ. **60**, 546–562 (2009). https://doi.org/10.1111/j.1477-9552.2009.00210.x
9. Tyllianakis, E., Martin-Ortega, J.: Agri-environmental schemes for biodiversity and environmental protection: how we are not yet "hitting the right keys." Land Use Policy **109**, 105620 (2021). https://doi.org/10.1016/J.LANDUSEPOL.2021.105620
10. Coyne, L., Kendall, H., Hansda, R., Reed, M.S., Williams, D.J.L.: Identifying economic and societal drivers of engagement in agri-environmental schemes for English dairy producers. Land Use Policy **101**, 105174 (2021). https://doi.org/10.1016/J.LANDUSEPOL.2020.105174
11. Defrancesco, E., Gatto, P., Mozzato, D.: To leave or not to leave? Understanding determinants of farmers' choices to remain in or abandon agri-environmental schemes. Land Use Policy **76**, 460–470 (2018). https://doi.org/10.1016/j.landusepol.2018.02.026

12. Espinosa-Goded, M., Barreiro-Hurlé, J., Ruto, E.: What do farmers want from agri-environmental scheme design? A choice experiment approach. J. Agric. Econ. **61**, 259–273 (2010). https://doi.org/10.1111/J.1477-9552.2010.00244.X

13. Capitanio, F., Adinolfi, F., Malorgio, G.: What explains farmers' participation in rural development policy in Italian southern region? An empirical analysis. New Medit. **10**, 19–24 (2011)

14. Mazorra, A.P.: Agri-environmental policy in Spain. The agenda of socio-political developments at the national, regional and local levels. J. Rural Stud. **17**, 81–97 (2001). https://doi.org/10.1016/S0743-0167(00)00028-0

15. Peerlings, J., Polman, N.: Farm choice between agri-environmental contracts in the European Union. J. Environ. Plan. Manag. **52**, 593–612 (2009). https://doi.org/10.1080/09640560902958131

16. Mora-herrera, D.Y., Huerta-barrientos, A.: A review of agent-based modeling for simulation of agricultural Una revisión de modelación basada en agentes para la simulación de sistemas agropecuarios. **88**, 103–110 (2021)

17. Schreinemachers, P., Berger, T.: An agent-based simulation model of human-environment interactions in agricultural systems. Environ. Model. Softw. **26**, 845–859 (2011). https://doi.org/10.1016/j.envsoft.2011.02.004

18. Brady, M., Sahrbacher, C., Kellermann, K., Happe, K.: An agent-based approach to modeling impacts of agricultural policy on land use, biodiversity and ecosystem services. Landsc. Ecol. **27**, 1363–1381 (2012). https://doi.org/10.1007/s10980-012-9787-3

19. Swinscoe, T.H.A., Knoeri, C., Fleskens, L., Barrett, J.: Agent-based modelling of agricultural water abstraction in response to climate change and policies: In East Anglia, UK (2014)

20. Matthews, R.B., Gilbert, N.G., Roach, A., Polhill, J.G., Gotts, N.M.: Agent-based land-use models: a review of applications. Landsc. Ecol. **22**, 1447–1459 (2007). https://doi.org/10.1007/s10980-007-9135-1

21. Kremmydas, D., Athanasiadis, I.N., Rozakis, S.: A review of agent based modeling for agricultural policy evaluation. Agric. Syst. **164**, 95–106 (2018). https://doi.org/10.1016/j.agsy.2018.03.010

22. Louhichi, K., Espinosa, M., Ciaian, P., Perni, A., Vosough Ahmadi, B., Colen, L., Gomez y Paloma, S.: The EU-wide individual farm model for common agricultural policy analysis (IFM-CAP v. 1). Economic impacts of CAP greening https://doi.org/10.2760/218047 (2018)

23. Huber, R., Bakker, M., Balmann, A., Berger, T., Bithell, M., Brown, C., Grêt-Regamey, A., Xiong, H., Le, Q.B., Mack, G., Meyfroidt, P., Millington, J., Müller, B., Polhill, J.G., Sun, Z., Seidl, R., Troost, C., Finger, R.: Representation of decision-making in European agricultural agent-based models. https://doi.org/10.1016/j.agsy.2018.09.007 (2018)

24. Ziv, G., Beckmann, M., Bullock, J., Cord, A., Delzeit, R., Domingo, C., Dreßler, G., Hagemann, N., Masó, J., Müller, B., Neteler, M., Sapundzhieva, A., Stoev, P., Stenning, J., Trajković, M., Václavík, T.: BESTMAP: behavioural, ecological and socio-economic tools for modelling agricultural policy. Res. Ideas Outcomes. **6**, 1–48 (2020). https://doi.org/10.3897/rio.6.e52052

25. Wittstock, F., Hötten, D., Biffi, S., Domingo-marimon, C., Bořivoj Šarapatka, M.B., Mesaroš, M.: Summaries of data, obstacles and challenges from interview campaigns (2020)

26. Grutters, J.P.C., Kessels, A.G.H., Dirksen, C.D., Van Helvoort-Postulart, D., Anteunis, L.J.C., Joore, M.A.: Willingness to accept versus willingness to pay in a discrete choice experiment. Value Health **11**, 1110–1119 (2008). https://doi.org/10.1111/j.1524-4733.2008.00340.x

27. Wilensky, U.: NetLogo. Northweste (1999)

28. Salecker, J., Sciaini, M., Meyer, K.M., Wiegand, K.: The nlrx r package: a next-generation framework for reproducible NetLogo model analyses (2019). https://doi.org/10.1111/2041-210X.13286

29. Morris, M.D.: Factorial sampling plans for preliminary computational experiments. Technometrics **33**, 161–174 (1991). https://doi.org/10.1177/001872086700900503

30. Campolongo, F., Cariboni, J., Saltelli, A.: An effective screening design for sensitivity analysis of large models. Environ. Model. Softw. **22**, 1509–1518 (2007). https://doi.org/10.1016/j.envsoft.2006.10.004
31. Saltelli, A., Tarantola, S., Campolongo, F., Ratto, M.: Sensitivity Analysis in Practice: A Guide to Assessing Scientific Models (Google eBook) (2004)

# Co-simulation of Socio-Technical Energy Systems: An Interdisciplinary Design Process

**Fabian Adelt, Matteo Barsanti, Sebastian Hoffmann, Debopama Sen Sarma, Jan Sören Schwarz, Ben Vermeulen, Tom Warendorf, Claudia Binder, Bert Droste-Franke, Sebastian Lehnhoff, Johanna Myrzik, Christian Rehtanz, and Johannes Weyer**

**Abstract** For exploration of future transition paths of the energy system and the complex challenges related to it, modeling components that are either a part of or connected to the energy system is primary. Here, co-simulation approaches facilitate integrated simulation scenarios by coupling simulation models developed in different programming languages, based on different modeling paradigms, and depicting various domains of the energy system (e.g., industry, households, or the electricity grid). However, co-simulation approaches exhibit a range of challenges and are thus under-exploited when investigating socio-technical transitions. We introduce a design and modeling process for an agent-based co-simulation framework, which aims to foster interdisciplinary collaboration considering multiple socio-technical elements of the energy system. This starts with building an information model for simulation planning and collecting inputs and outputs of different models. Finally, a modularization approach defines simulation sub-scenarios to simplify modeling interdependencies. Additionally, we present two exemplary scenarios: (i) the impact of households' energy-related behavior on power grid stability and (ii) the co-evolutionary supply and demand dynamics of energy storage technologies in the industrial sector.

F. Adelt · S. Hoffmann (✉) · J. Weyer
Technology Studies Group, Faculty of Social Sciences, TU Dortmund University, Dortmund, Germany
e-mail: sebastian3.hoffmann@tu-dortmund.de

M. Barsanti · C. Binder
School of Architecture, Civil and Environmental Engineering, EPFL, Ecublens VD, Switzerland

D. S. Sarma · C. Rehtanz
Institute of Energy Systems, Energy Efficiency and Energy Economics, TU Dortmund University, Dortmund, Germany

J. S. Schwarz · S. Lehnhoff
Department of Computing Science, University of Oldenburg, Oldenburg, Germany

B. Vermeulen · B. Droste-Franke
Institut für qualifizierende Innovationsforschung und -beratung GmbH (IQIB), Bad Neuenahr-Ahrweiler, Germany

T. Warendorf · J. Myrzik
Institute of Automation, University of Bremen, Bremen, Germany

477

## 1  Introduction

The energy system represents a complex, socio-technical system whose transition to
low-carbon energy involves changes within several domains and dimensions, such
as the demand side (consumer practices), the supply side (infrastructure, technologies, business models), the intermediate layers (storage, transmission, distribution,
trading), and institutional dimensions (policies, regulations). In this regard, energy
models are *"powerful tools for a systematic, quantitative and forward-looking analysis"* [8, p. 162] that might help to inform policy-making [18] to approach these
complex challenges. Typically, this involves evaluating future transition pathways
or scenarios that seek to achieve predefined, sustainability-related targets (i.e., environmental or economic impact, social acceptance, and security of supply) [11]. The
suitability of different energy modeling approaches has already been investigated and
compared in recent studies [8, 12], which highlight how each paradigm represents
a particular perspective on energy modeling—each associated with its respective
prospects and challenges, and no specific method being superior to another [12].

The majority of existing energy system simulation models (cf. reviews in [3, 4])
either address particular problems and objectives (e.g., generation capacity, renewables mix, grid optimization) or represent a specific part of the energy sector (e.g.,
residential customers, wind power). Consequently, domains not in focus are often
included as boundary conditions or based on simplifying assumptions. In this context,
researchers note the need for more realistic models that integrate multiple elements
(e.g., political, economic, technological) as well as their interactions and co-evolution
[12]. Some authors particularly emphasize the lack of adequate incorporation of
social elements in energy models [6] and voice the need for more early interdisciplinary collaboration [24]. With regard to the latter, authors have already proposed
the application of common socio-technical concepts and terminologies to foster the
exchange in interdisciplinary teams [7], the sharing of tools, data, and strategies
to increase transparency, accessibility, and re-usability of models [15, 20], or the
utilization of disciplinary expertise to circumvent respective drawbacks [8].

We follow up on this discussion and intend to raise awareness of *co-simulation*,
which aims to couple diverse models in integrated simulation scenarios and which has
not yet received adequate attention in modeling socio-technical transitions. According to a keyword analysis on Scopus by [23], related literature has been steadily
growing over the last decades. However, the use of energy system co-simulation is
still more common in technical and mathematical disciplines than in computational
social sciences, indicating a need for further research with respect to the inclusion
of human and social behavior models [26].

This paper presents potentials and challenges of co-simulation as a tool for modeling socio-technical energy transitions. Based on the authors' interdisciplinary project

MoMeEnT (*"Modelling the socio-technical multi-level architecture of the energy system and its transformation"*), a design process is introduced to support modelling activities. This process starts on a content level by collecting data in an information model and then proceeds step-wise towards setting up co-simulation scenarios based on general research aims and built-in hypotheses.

We introduce the general idea of co-simulation in Sect. 2, address our own co-simulation approach in Sect. 3, and provide exemplary scenarios in Sect. 4. In Sect. 5, we conclude with a discussion of our proposed process and mention crucial aspects that need consideration when adopting it.

## 2 Background: Co-simulation Approaches and Concepts

Generally, co-simulation encompasses a set of coupled simulators, each implementing a model intended to represent real-world entities, phenomena, or processes (e.g., the power distribution grid, industrial processes, or households) [17]. Collectively, these simulators form a complex system. Dynamic interconnections and interdependencies between simulators are established through data interfaces, i.e., the definition of data flows. Thus, the simulators can operate simultaneously, are only loosely coupled, and retain individual time step sizes.

A variety of methods and concepts for co-simulation exists [5, 16], but often the various simulators are not linked directly within a co-simulation. Usually, a central "master algorithm" or co-simulation framework coordinates the setup and initialization of the simulators, the order of execution, the time step synchronization, and the exchange of data [17]. According to [17, p. 40], the actual "art" of co-simulation lies within this orchestration of simulators.

From the authors' perspective, using co-simulation for modeling socio-technical transitions brings the following benefits: Firstly, instead of a monolithic simulation of a single domain, co-simulation inherently supports holistic simulation that integrates multiple domains with sub-system interactions. Secondly, it allows reusing established, well-suited, and potentially well-developed tools for specific domains instead of developing new integrated models for every use case. Thirdly, the scenario definition based on connecting simulator in- and outputs is easier than integrated modeling approaches. Finally, in contrast to an integrated approach, simulators may remain black boxes. Consequently, modelers from different disciplines do not need to know the specifics of the other simulators.

Nonetheless, co-simulation also has challenges compared to other simulation approaches, from which we see the following three as the most important. **Issue 1**: Coupling diverse models can still be complex and demanding if there are many model interdependencies. **Issue 2**: The complexity of multiple coupled models can also complicate the validation of scenarios, might decrease the robustness of a simulation, and can reduce the performance. **Issue 3**: A core task of co-simulation frameworks is to harmonize the time steps of all models. Usually, transition processes are long-term scenarios where the models have diverging time scales which need specific handling. In the next section, we address these issues using our proposed approach.

# 3 Proposed Co-simulation Design Process and Workflow

The design of an interdisciplinary co-simulation scenario can be challenging because many data flows and interdependencies of the models need consideration and expected simulation outputs need to be defined (Issue 1). Even if the modelers from different domains can work independently on their models, a collaboration between multiple domain experts is crucial for this design process. For assistance, a previously developed approach, based on an information model for simulation planning [21, 22], was integrated with a new method to define general aims and hypotheses for the simulation. The overall methodological framework is depicted in Fig. 1, summarizing the objectives of each process step (O) and the supporting tools (T).
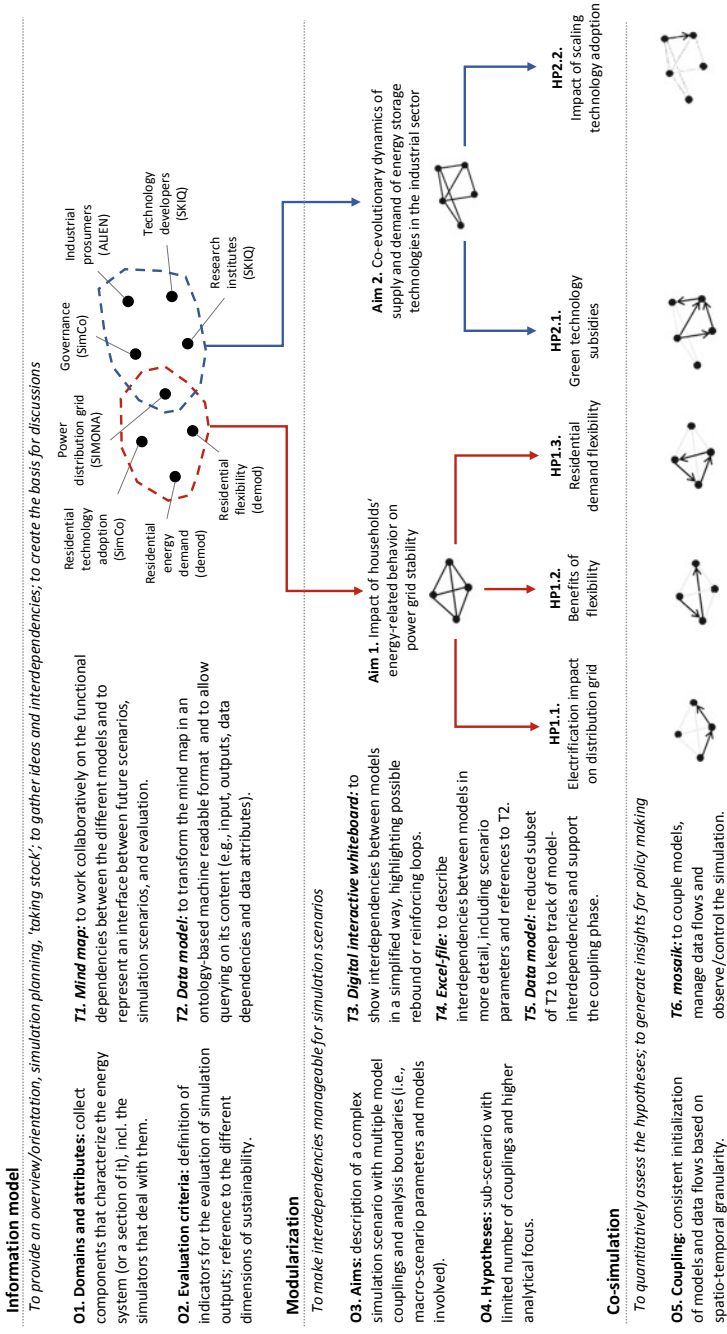
The information model provides an ontological structure to describe the data flows between models [22]. We started the collaborative design process with model collection and data flow modeling in the form of a mind map, later converted to a structure that allows automated content export in the ontological data model, enabling further processing and information extraction to support the scenario planning process. Nevertheless, dialogue on the common understanding of terms and concepts is still necessary, e.g., to avoid conceptual redundancies between models, but can be supported by existing domain ontologies.

Although co-simulation allows the integration of heterogeneous models into a holistic socio-technical system, this representation has a cost (Issue 2): It requires validation on both the micro-level of individual sub-models and the macro-level of the coupled system or scenario. According to different authors, judging the latter form of validity represents a considerable task in co-simulation, and there is still a need for further research regarding the representation and enforcement of model validity assumptions [16, 17, 23].

To overcome such complexity and validation issues on the macro-level, we decided to include an intermediate 'meso-level validation' by modularizing the full scenario (i.e., all eight simulators in Fig. 1) into smaller sub-scenarios. Model dependencies within a simulation scenario thus become more manageable and allow for focused investigation and discussion of specific model interactions, ultimately helping to enhance content validity.

As a first step, we discussed critical challenges, which require a limited number of sub-system models (i.e., a smaller group of simulators) coupled within a simulation scenario. Such challenges—or "bottlenecks"—arise, for example, from current inefficiencies, unused potentials, or path dependencies that impede or even prevent low-carbon transitions towards sustainability [13]. After interdisciplinary discussions, we formulated various 'general aims' (i.e., research questions) to be addressed by different simulator sub-sets. In the early phase of the modeling process, these general aims provided framing for the further specification of more manageable co-simulation sub-scenarios and their analytical boundaries (e.g., simulated time horizon). Two exemplary general aims are discussed in Sects. 4.1 and 4.2.

Since the general aims still describe complex simulation scenarios, we transformed them into more specific hypotheses: These represent sub-scenarios that deal

**Fig. 1** Framework of the project methodology. Note: (O) objective; (T) supporting tool; (HP) hypothesis

with the assumed interdependencies between in- and output variables of a limited number of simulators. Consequently, hypotheses can be linked to the information model to show which data flows are needed for their implementation.

Nevertheless, the relatively small number of design steps should not obscure the fact that forming such general aims and hypotheses is a lengthy process of discourse and coordination between the involved modelers. The "black-box" character of co-simulations (cf. [23]) benefits modelers to collaborate without intricate knowledge of models outside their discipline. However, an exchange of (implicit) assumptions within and between the individual simulators is necessary to increase the consistency of co-simulation scenarios. On the one hand, a collaborative whiteboard platform (i.e., Miro) allowed us to visualize interdependencies between simulators in a simplified way, highlighting possible overlaps or reinforcing loops in the context of a general aim. On the other hand, information exported from the mind map (i.e., the various model components) was collected in a spreadsheet (i.e., Microsoft Excel) and used as 'building blocks' for deriving hypotheses and specifying the implemented scenarios within the co-simulation framework *mosaik*.

## Co-simulation framework: *mosaik*

For setting up co-simulations many different frameworks exist, as compared in [27]. The framework *mosaik*[1] was selected here because of its focus on usability and flexibility, as described in [25], and because it is available as open-source software. The core of *mosaik* is Python-based and provides a *component API* to connect models from diverse programming languages and integrate them into scenarios with a *scenario API*.

The coupling of diverse models in a co-simulation increases complexity, hampers the validation and robustness of simulation, and raises performance issues (Issue 2). For a robust simulation, we propose the observer/controller architecture described in [19], which observes parts of the simulation for critical states and reacts under specific conditions or takes control. Performance, albeit not the primary focus of *mosaik*, can be addressed by distributing computation to multiple computers and optimizing the temporal aspects of a scenario, e.g., by calculating a set of representative days instead of doing a full-year simulation.

Issue 3 is generally addressed by co-simulation frameworks, as it is one of their primary tasks to orchestrate the different simulators by organizing their data flows and execution times. For substantial differences in time scale, a new *mosaik-aggregator* component has been developed. It can be placed between two simulators and aggregates outputs of one simulator over time and provides these aggregated results as input to the other simulator. Thus, simulators do not need to adapt to different time scales and reduce their universality.

---

[1] https://mosaik.offis.de.

## 4   Example Scenarios

Here, to provide a concrete application of the process proposed above, we present two exemplary scenarios, describe the simulators involved and discuss the challenges encountered in their coupling.

### *4.1   Scenario 1: Electrification of Residential End-Users*

The first scenario investigates the effects of electrification and flexibility increase in residential energy demand, focusing on space heating with heat pumps, electric vehicle charging, and distributed generation from photovoltaic (PV) systems. Modularization generates three interrelated hypotheses: (HP1.1) increased investments in electrical technologies lead to an increase in grid volatility; (HP1.2) Demand-Side Management (DSM) program benefits (i.e., lower energy costs and increased self-consumption) foster the adoption of low carbon technologies; and (HP1.3) consumers' active engagement in DSM programs alleviates congestion risk due to low-carbon technology adoption. It is noteworthy that this list of hypotheses is not exhaustive. However, it represents a good starting point for understanding the behavior of coupled models and exploring emerging dynamics.

**Simulators**. This scenario contains two types of simulators. One deals with residential technology adoption and energy consumption patterns, and the other represents a distribution grid to assess the technical aspects of the hypotheses.

*SimCo-Energy*[2] aims to model households' long-term adoption decisions in various energy-related domains, specifically PV systems, household appliances, heating systems, and car-based mobility. To capture the varying complexity and rationality of investment decisions (e.g., habitual or reflecting-calculating) as well as social influences (e.g., conformity and trust), households' decision rules follow the Consumat approach [10] and its ABMS-implementation by [14]: Households evaluate the fulfillment of their social and personal needs by comparing the perceived performances of available technologies with regard to multiple evaluation criteria (e.g., financial consequences, environmental impact, or effort). An online survey is used to collect empirical data for calibrating the household agents (e.g., concerning preferences).

Residential energy demand is simulated by *demod*,[3] a Python library that provides a set of modules to assemble Germany-based domestic energy demand models in a customizable, fully documented, and open-source manner. It implements an activity-based approach to energy demand modeling. From time-use statistics, activity profiles of individual household members are generated stochastically with a time resolution of 10 min. Household activity profiles, technical characteristics, and sta-

---

[2] SimCo-Energy is an adaption of the general-purpose ABMS framework SimCo ("Simulation of the governance of complex systems") [1].

[3] https://demod.readthedocs.io/en/latest/.

tistical data of different appliances are then used as input to simulate switch-on/off events and the daily electrical load profile with a time resolution of 1 min. An application of this model in a co-simulation scenario is provided in [2].

*SIMONA*[4] is an agent-based simulation framework representing modern distribution networks and aids in grid planning and operation. The software is written in SCALA and uses the AKKA agent framework. The primary feature of *SIMONA* is to provide distribution grid forecasts using different load profiles. It efficiently solves a distributed power flow across multiple voltage levels connected by transformers [9]. The results obtained from the power flow calculations provide valuable insights into how the applied load affects grid parameters such as nodal voltages and line loadings. Consequently, one can gauge the congestion level and the risk to the security of supply in a distribution grid.

**Coupling**. In the initialization phase, the input data of the simulators (e.g., population size and grid characteristics) must generate a plausible simulation scenario for each. For example, undersizing the grid relative to population size or PV or EV penetration makes the simulation impractical since *SIMONA* will generate solutions incompatible with grid constraints. Second, during the simulation of two or more behaviors/dynamics belonging to the same agent, the characterizations of this agent in the different simulators must be consistent. For example, although *SimCo* and *demod* treat two distinct consumer behaviors, they share some variables related to their profiling, such as socio-psychological factors, which requires care at the parametrization stage, data collection planning, and processing.

In the simulation phase, it is necessary to couple the inputs and outputs of the different simulators. This requires adapting the data format and additional processing to handle the models' different spatial and temporal scales. Indeed, *demod* can generate both appliance-level and household-level electricity profiles with 1 min resolution. This data must be aggregated (via *mosaik-aggregator*) every 30 min to provide the average household consumption needed by *SIMONA* to solve the network power balance equations and every month to provide the monthly energy consumption per device to *SimCo* agents.

Ideally, the coupling results should show the importance of the co-evolutionary dynamics between low-carbon technology adoption behavior, device usage, and distribution grid operation in the transition towards a decarbonized, electrified, demand flexibility-based energy system. Therefore, a better understanding and quantification of consumer responses (e.g., to blackout risks due to grid congestion) will help form new energy interventions and policies.

## *4.2 Scenario 2: Electrification of Industrial Prosumers*

The second scenario aims to study the co-evolutionary dynamics of supply and demand of energy storage technologies and the existence of so-called waiting games

---

[4] https://simona.ie3.e-technik.tu-dortmund.de.

and volatility in market tipping behavior. This combination of models studies the impact of government interventions such as taxation and subsidies on the development and adoption of storage technologies, which supports energy policy recommendations. Thus, we formulate two hypotheses: (HP2.1) supporting industrial research and development of green technology drives its adoption and thereby the energy transition; and (HP2.2) an increase in industrial adoption of renewable energy technology amplifies its development by increasing scale and lowering unit costs.

**Simulators**. This scenario uses three simulators, the active local industry energy network *(ALIEN)*, an energy technology innovation process model *(SKIQ)*, and a governance module (as part of *SimCo*).

The goal of *ALIEN* is to represent industrial companies and their energy usage as well as their energy management agents (EMA). They actively cost optimize their energy consumption and production by managing their self-load, solar power production, and flexibility from storage.

In brief, *SKIQ* has a population of firms engaged in developing technologies reproducing empirically-grounded learning curves in key performance parameters and price. In this project, *SKIQ* simulates firms involved with developing and selling energy storage technologies to customers (represented by *ALIEN*), buying these technologies to reduce energy costs.

*SimCo's* governance module entails a set of possible interventions (like technology bans or specific subsidies), which activate based on boundary conditions like overall emissions of industry, market penetration of certain old or new technologies, or percentage of renewable generation in the system.

**Coupling**. The co-simulation framework *mosaik* initializes the simulators by specifying the size of the population of buyers (*ALIEN*) and suppliers (*SKIQ*) and their respective endowments (available capital, energy consumption/generation profile, battery technology specialization, etc.) as well as a set of possible government interventions (*SimCo*) and their activation conditions. During the simulation, *mosaik* functions as a communication channel for the sets of technologies offered (from *SKIQ* to *ALIEN*) and information on which are purchased and required (from *ALIEN* to *SKIQ* and *SimCo*).

During the process of coupling, several conceptual issues came to light. *ALIEN* is a 'technical' decision support tool that models the adoption of energy storage technology by a *single* industrial firm in extensive detail. In contrast, *SKIQ* is a 'social' simulation model in which a *population* of technology developers are making and selling energy storage technologies on the market to a *range* of industrial firms. To couple both simulators, a simplified population of industrial firms is needed to connect to *SKIQ*. This population simplifies the industrial firms' energy management by excluding uncertainty to reduce computational expenses. However, due to the short-term focus of the industrial agents, *ALIEN* lacks forward-looking decisions and the formulation of technical requirements. Consequently, technology developers in *SKIQ* make their technology development decisions based on actual purchases and not on needs, thereby not on potential future purchases. So, with information shortage on (yet unmet) requirements from the industry, decisions on the technology

development direction by agents in *SKIQ* are reactive rather than proactive. Further-more, *ALIEN* and *SKIQ* have different time resolutions, which must be aggregated by *mosaik*.

The *ALIEN-SKIQ* simulation results reveal segments of the parameter space in which desirable technologies have a strong tendency to become dominant at an acceptable rate. Depending on the scenario, government intervention is needed to tip the market toward desired technologies and accelerate their development and adoption.

## 5    Discussion and Conclusion

In this paper, we presented the potentials and challenges of co-simulation as a tool for modeling energy transitions and exploring co-evolutionary dynamics between different socio-technical elements of the energy system. As interdependencies increase, coupling becomes more complex, both in terms of conceptualization (models) and implementation (simulators). Therefore, we proposed a design process to support modelling activities and simulation planning, which is based on an information model. Using mind maps with an ontological data model is an easy-to-use approach to visualize the simulation infrastructure and simultaneously perform targeted queries about the type of data flows between simulators. Thus, it eases collaboration and communication among multi-disciplinary researchers without extensive familiarization with specific modeling tools.

In addition, to address validation challenges at the micro/macro-level (at the level of individual sub-models and coupled systems or scenarios, respectively), we supplemented the information model with a problem modularization strategy based on *general aims* and *hypotheses*. Accordingly, instead of using an integrated, complete co-simulation scenario, we divide the overarching research question into separate sub-scenarios that use a subset of simulators. Each sub-scenario or *general aim* seeks to study a possible "bottleneck" of the energy system transition where a limited number of sub-system model couplings is required. While structuring the problem into *general aims* simplifies the definition of analysis boundaries and scenario parameters and reduces the number of coupled simulators, validating and understanding the co-simulation remains a complex task. Thus, we further modularized the problem by extracting a set of *hypotheses* for each *general aim*, studying dynamics in isolation without complications arising from co-evolution with other simulators. It is important to note that this also simplifies cross-validation with (limited) formal models, comparison with historical data, and verification of stylized facts.

Although this work proposes a methodological approach to overcome some challenges in co-simulation design, we cannot overlook the following limitations: Firstly, while co-simulation facilitates the distribution of modeling activities over domain experts and programming in different languages, modelers still require a conceptual understanding of distinct simulators beyond the technical specifications of the interfaces. Secondly, in repurposing simulators, the adequacy of internal simulator

data and model assumptions needs to be reassessed, especially when time horizon or spatial setting changes between simulator deployments, the stationary of certain variables may be implausible given underlying trends, agent heuristics and objectives may change, the time resolution of simulation is impractical, or specific value ranges are not appropriate.

To conclude, we believe it is important to emphasize that the co-simulation approach, besides its original function of coupling and orchestrating independent simulators, provides a concrete common basis for the critical analysis of socio-technical energy system models through the joint exploration and unambiguous specification of research questions, simulator purposes and operational assumptions by various domain experts. Indeed, only through an iterative process in which the interfaces, purposes and operational specifications of subsystem models are informed in a multidisciplinary manner can one aspire to model the complexity and nonlinearity of the interconnected actor-dense socio-technical energy system.

# References

1. Adelt, F., Weyer, J., Hoffmann, S., Ihrig, A.: Simulation of the governance of complex systems (SimCo): basic concepts and experiments on urban transportation. J. Artif. Soc. Soc. Simul. **21** (2018)
2. Barsanti, M., Schwarz, J.S., Gérard Constantin, L.G., Kasturi, P., Binder, C.R., Lehnhoff, S.: Socio-technical modeling of smart energy systems: a co-simulation design for domestic energy demand. Energy Inform. **4**(3), 12 (2021)
3. Connolly, D., Lund, H., Mathiesen, B.V., Leahy, M.: A review of computer tools for analysing the integration of renewable energy into various energy systems. Appl. Energy **87**(4), 1059–1082 (2010)
4. Fattahi, A., Sijm, J., Faaij, A.: A systemic approach to analyze integrated energy system modeling tools: a review of national models. Renew. Sustain. Energy Rev. **133**, 110195 (2020)
5. Hafner, I., Popper, N.: On the terminology and structuring of co-simulation methods. In: Zimmer, D., Bachmann, B. (eds.) Proceedings of the 8th International Workshop on Equation-Based Object-Oriented Modeling Languages and Tools, pp. 67–76. ACM, New York, NY, USA (2017)
6. Hansen, P., Liu, X., Morrison, G.M.: Agent-based modelling and socio-technical energy transitions: a systematic literature review. Energy Res. Soc. Sci. **49**, 41–52 (2019)
7. Hinker, J., Hemkendreis, C., Drewing, E., März, S., Rodríguez, D.I.H., Myrzik, J.M.A.: A novel conceptual model facilitating the derivation of agent-based models for analyzing socio-technical optimality gaps in the energy domain. Energy **137**, 1219–1230 (2017)
8. Hirt, L.F., Schell, G., Sahakian, M., Trutnevyte, E.: A review of linking models and socio-technical transitions theories for energy and climate solutions. Environ. Innov. Societal Transit. **35**, 162–179 (2020)
9. Hiry, J., Kittl, C., Sarma, D.S., Oberließen, T., Rehtanz, C.: Multi-voltage level distributed backward-forward sweep power flow algorithm in an agent-based discrete-event simulation framework. Electr. Power Syst. Res. **213**, 108365 (2022)

10. Jager, W., Janssen, M.: The need for and development of behaviourally realistic agents. In: Goos, G., Hartmanis, J., van Leeuwen, J., Sichman, J.S., Bousquet, F., Davidsson, P. (eds.) Multi-Agent-Based Simul. II, vol. 2581, pp. 36–49. Springer, Berlin Heidelberg (2003)

11. Jewell, J., Cherp, A., Riahi, K.: Energy security under de-carbonization scenarios: an assessment framework and evaluation under different technology and policy choices. Energy Policy **65**, 743–760 (2014)

12. Köhler, J., de Haan, F., Holtz, G., Kubeczko, K., Moallemi, E., Papachristos, G., Chappin, E.: Modelling sustainability transitions: an assessment of approaches and challenges. J. Artif. Soc. Soc. Simul. **21**(1) (2018)

13. Loock, M.: Unlocking the value of digitalization for the European energy transition: a typology of innovative business models. Energy Res. Soc. Sci. **69**, 101740 (2020)

14. Moglia, M., Podkalicka, A., McGregor, J.: An agent-based model of residential energy efficiency adoption. J. Artif. Soc. Soc. Simul. **21** (2018)

15. Morrison, R.: Energy system modeling: public transparency, scientific reproducibility, and open development. Energy Strategy Rev. **20**, 49–63 (2018)

16. van Nguyen, Besanger, Y., Tran, Q., Nguyen, T.: On conceptual structuration and coupling methods of co-simulation frameworks in cyber-physical energy system validation. Energies **10**(12), 1977 (2017)

17. Palensky, P., van der Meer, A.A., Lopez, C.D., Joseph, A., Pan, K.: Cosimulation of intelligent power systems: fundamentals, software architecture, numerics, and coupling. IEEE Industr. Electron. Mag. **11**(1), 34–50 (2017)

18. Papachristos, G.: Towards multi-system sociotechnical transitions: why simulate. Technol. Anal. Strat. Manage. **26**, 1037–1055 (2014)

19. Richter, U., Mnif, M., Branke, J., Müller-Schloer, C., Schmeck, H.: Towards a generic observer/controller architecture for organic computing. GI Jahrestagung **93**, 112–119 (2006)

20. Schulze, J., Müller, B., Groeneveld, J., Grimm, V.: Agent-based modelling of social-ecological systems: achievements, challenges, and a way forward. J. Artif. Soc. Soc. Simul. **20** (2017)

21. Schwarz, J.S., Elshinawy, R., Ramírez Acosta, R.P., Lehnhoff, S.: Ontological integration of semantics and domain knowledge in hardware and software co-simulation of the smart grid. In: Fred, A., Salgado, A., Aveiro, D., Dietz, J., Bernardino, J., Filipe, J. (eds.) Knowledge Discovery, Knowledge Engineering and Knowledge Management, vol. 1297, pp. 283–301. Springer, Cham (2020)

22. Schwarz, J.S., Witt, T., Nieße, A., Geldermann, J., Lehnhoff, S., Sonnenschein, M.: Towards an integrated development and sustainability evaluation of energy scenarios assisted by automated information exchange. In: Donnellan, B., Klein, C., Helfert, M., Gusikhin, O., Pascoal, A. (eds.) Smart Cities, Green Technologies, and Intelligent Transport Systems, vol. 921, pp. 3–26. Springer, Cham (2019)

23. Schweiger, G., Gomes, C., Engel, G., Hafner, I., Schoeggl, J., Posch, A., Nouidui, T.: An empirical survey on co-simulation: promising standards, challenges and research needs. Simul. Modell. Pract. Theor. **95**, 148–163 (2019)

24. Süsser, D., Gaschnig, H., Ceglarz, A., Stavrakas, V., Flamos, A., Lilliestam, J.: Better suited or just more complex? on the fit between user needs and modeller-driven improvements of energy system models. Energy **239**, 121909 (2022)

25. Steinbrink, C., Blank-Babazadeh, M., El-Ama, A., Holly, S., Lüers, B., Nebel-Wenner, M., Ramírez Acosta, R., Raub, T., Schwarz, J.S., Stark, S., Nieße, A., Lehnhoff, S.: CPES testing with Mosaik: co-simulation planning, execution and analysis. Appl. Sci. **9**(5), 923 (2019)

26. Steinbrink, C., Schlogl, F., Babazadeh, D., Lehnhoff, S., Rohjans, S., Narayan, A.: Future perspectives of co-simulation in the smart grid domain. In: 2018 IEEE International Energy Conference, ENERGYCON 2018, pp. 1–6 (2018)

27. Vogt, M., Marten, F., Braun, M.: A survey and statistical analysis of smart grid co-simulations. Appl. Energy **222**, 67–78 (2018)

# Dynamics of Individual Investments in Heating Technology

**Sascha Holzhauer, Friedrich Krebs, and Lukas Jansen**

**Abstract** The transition of heat provision in the urban building stock towards climate neutral sources poses a major challenge to German cities. The underlying actor structure is complex and interlinked. Municipalities set regulatory boundary conditions and decide on infrastructure investments like district heating networks. Necessary investments on the premises of house owners are not only inhibited by unavailability of capital but also by a lack of technical knowledge and ultimately by capacity shortages of installation companies. In the paper, we outline the agent-based model being developed in the course of a new research project aiming to support local heat transitions by socio-technical modelling and simulation. We aim to represent the investment dynamics evolving from interactions of building owners with a broader set of stakeholders, namely energy consultants whose knowledge and thus recommendations shape the set of investment options, and craftspeople such as plumbers whose experience has an impact of building owners' decisions. We outline the integration of the agent-based model with a model of the local energy system to account for feedbacks between the heating infrastructure and investment decisions of building owners. Furthermore, we discuss our approach to auto-parameterise intervention measures to achieve required rates of investments.

**Keywords** Heat transition · Local energy system analysis · Agent-based modelling · Recommendation networks · Theory of planned behaviour

S. Holzhauer (✉) · F. Krebs · L. Jansen
Section Integrated Energy Systems, University of Kassel, Wilhelmshöher Allee 73, 34121 Kassel, Germany
e-mail: Sascha.Holzhauer@uni-kassel.de

F. Krebs
Fraunhofer Institute for Energy Economics and Energy System Technology, Joseph-Beuys-Straße 8, 34117 Kassel, Germany

# 1 Motivation

Achieving German climate protection targets by 2030 requires stronger incentives for energetic renovation as well as a significant increase of renewable energy for heat provision in buildings. The latter can be achieved by switching to heat pumps as well as extensions of local and district heating grids fed by a high proportion of renewable energies. The need to act on municipality level has been recognized by the German government [1]. The goal of the project WAERMER[1] is to apply a methodical integration of energy system optimisation on a district level and agent-based modelling of individual investment behaviour to help identify the requirements for a successful transition of heat provision in the urban building stock towards $CO_2$-neutrality.

Often, municipalities plan and realise energy systems based on technical and economic aspects. However, it has been demonstrated that social impact factors need to be considered to achieve a successful energy transition [2]. In particular, the heat transition in the building sector requires mobilisation of capital flows from individual house owners, who face diverse internal and external barriers inhibiting investments in innovative heating technology. For instance, they lack equity capital or are not aware of attractive funding to guarantee profitable investments. In addition, they may avoid complex planning and building permission processes. Further, lack of knowledge of the renovation process, benefits and alternative options such as stepwise transitions or joint neighbourhood heating solutions hinder decision makers to renovate and exchange heating technology. On the other hand, motivation to invest in new heating technology or to renovate is often based on individual preferences such as comfort, environmental awareness, energy cost savings, and health issues [3].

Agent-based modelling (ABM) enables investigations about the interplay between external conditions (regulatory framework, funding opportunities, information provision) and investment dynamics on the part of individuals. Results have the potential to provide stakeholders and decision makers in policy, economy and civil society with evaluations and a-priori impact assessments of technical, regulatory, and business policy options [4].

Often, single measures do not achieve a significant increase in adoption rates, but combinations of fiscal policies and subsidies do [5], and also informal measures such as information campaigns are successful. A reason can be seen in the multi-stage character of adoption, when initial interest is succeeded by gaining knowledge, followed by the planning phase and the actual decision, before the new technology is finally implemented [6]. Each of these stages needs to be triggered and passed through to end up with the adoption of a new heating technology. ABM is able to explore the interplay of interventions of various kinds, which has been found crucial for their success to foster building renovation measures [7].
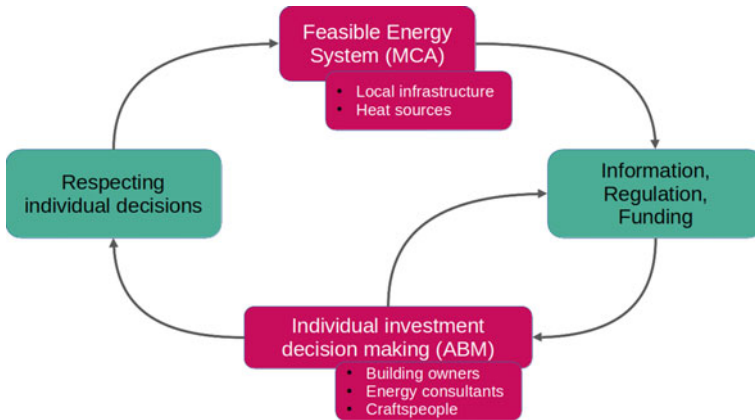
---

## 2 Modelling Approach

The focus of the modelling effort is on the identification and mitigation of obstacles for individual investors under consideration of relevant stakeholders such as municipality, public utility companies, craftspeople such as plumbers, and energy consultants. To integrate technical and socio-economic aspects, the agent-based model is part of a common simulation cycle with a planning tool for local heat provision (Fig. 1). The agent-based model outputs time series of investments per heating technology and area. Detailed knowledge of the local building stock such as year of construction, age of building parts and heating system, area, and insulation standard are required to trigger investments, identify suitable heating options and determine eligibility of funding. Therefore, a combined approach of synthetic reconstruction (SR) and combinatorial optimisation (CO) is applied to generate a realistic local building stock from German census data.

### 2.1 Integration of ABM with Energy System Modelling

The iterative coupling of the agent-based model with the local energy system modelling addresses the feedbacks between infrastructural decisions such as implementing district heating or allowing earth drilling to access heat sources, and individual decision-making based on heterogeneous investment preferences about heating technologies. Existing infrastructure influences the set of heating options

**Fig. 1** Based on a multi criteria analysis (MCA), the local energy system modelling tool suggests a feasible and optimal heat supply solution. Subsequently, the agent-based model simulates building owners' decision-making as well as influences from energy consultants and craftspeople. Interventions can be explored to incentivise individuals towards the optimal solution. In case these are not sufficient, the MCA needs to respect individual decisions to come up with an alternative solution

**Fig. 2** Building owners technology choices and local infrastructure should align for an effective and cost-efficient operation: local or district heating with a high proportion of connected buildings, or individual heat supply solutions without investments in common infrastructure
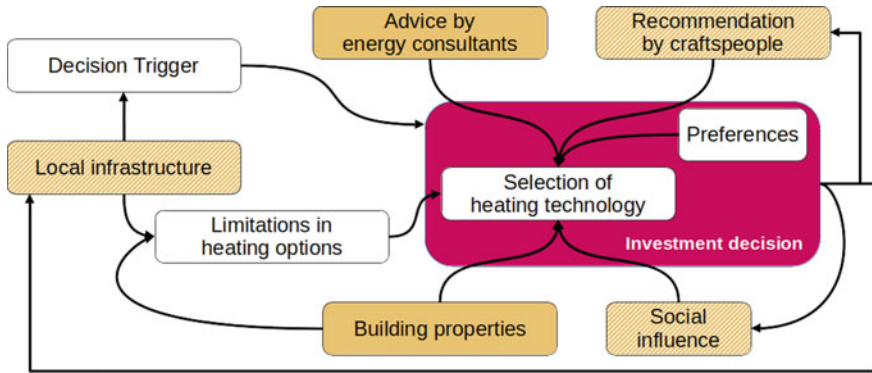
available to investing house owners, while individual investment decisions regarding the buildings connected to the infrastructure influence the effectiveness of infrastructure measures. By coupling the models, inconsistencies between the decisions made by municipalities and building owners can be detected and addressed (Fig. 2). For instance, the efficient operation of district heating involves a certain proportion of buildings connected to the grid. The local energy system model may determine the required number of connections, and by the ABM it can be explored which measures and combinations are necessary to achieve that proportion, and how these need to be parameterised.

## 2.2 Modelling Investment Decisions

In Germany, 85% of the residential building stock is privately owned [8]. The investment decisions of private building owners are shaped by numerous factors. Figure 3 shows an influence diagram that reflects our initial understanding on the multi-stage decision-making of private building owners.
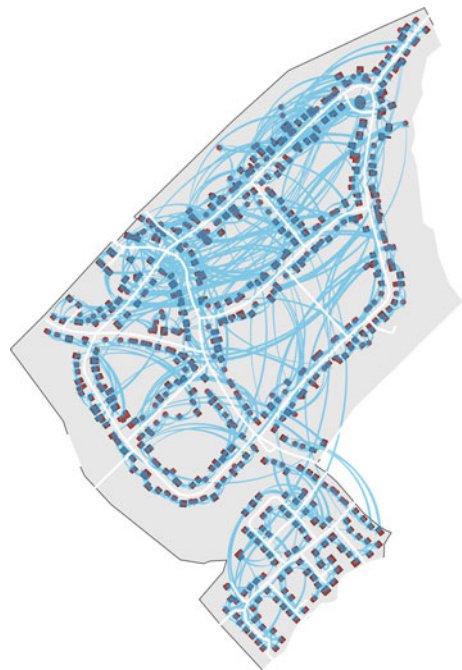
Factors influencing building owners' selection of a heating technology comprise the influence of plumbers, energy consultants, and their additional social network [9, 10], as well as individual preferences for cost reduction, environment-friendliness, low planning effort, low installation effort, and comfort. In addition, the model represents the owners' expertise regarding technologies and awareness of funding. Concerning energy consultants and craftspeople, their expertise and preferences towards specific heating technologies, open-mindedness, and training habits are deemed relevant to building owners' decision-making.

Furthermore, the model incorporates feedbacks between installations of certain kinds, the local infrastructure, and gained experience of and recommendations for specific technologies via social networks of energy consultants, craftspeople and peers. A spatially explicit and dynamic social network modelling provides the structure of social influence relations (Fig. 4). This allows the exploration of impacts on the frequency, intensity and content of knowledge exchange and recommendations.

**Fig. 3** Influence diagram of the building owners' decision making about investment in heating technology. Amber boxes mark external inputs to the agent-based model; hatched boxes are also subject to internal dynamics. The local infrastructure enables options (e.g., district heating), according interventions trigger decisions but also consider building owners' selections. Craftspeople's recommend depending on their experience, but also react to owners' demand and train accordingly. Social influence by peers is subject to their decisions and experiences

**Fig. 4** One of the study areas in Kiel/Germany. Buildings are shown in red, and blue arrows represent social network links of information exchange and influence between building owners

The implementation will be based on the LARA framework [11]. Private building owners' decision making is modelled as a combination of heuristics and an operationalisation of the Theory of Planned Behaviour [12], which has been proven to be suitable to incorporate empirical data into agent-based modelling [13, 14]. Figure 3 depicts some of the constructs of the theory: For instance, individual preferences guiding the investment decision may be understood as the subjective attitude towards the behaviour, dynamic social influence constructs a subjective norm, and local infrastructure and building properties are proxies for behavioural control beliefs of the investor.
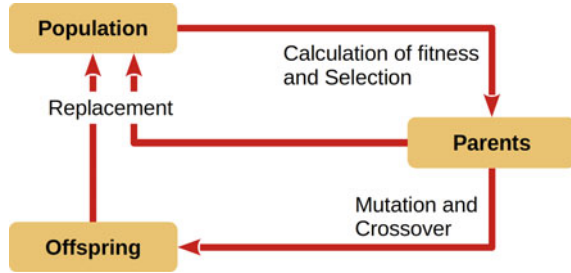
### 2.2.1  Auto-Parameterisation

In the project, we seek to find applicable setups for a successful heat transition at the community level. This often requires appropriate investment decision by a sufficient share of building owners. E.g., profitable public investments in district heating require the number of connected buildings to exceed a certain threshold. In case the baseline scenario finds that the number of building owners willing to connect to district heating falls short of this threshold, measures which incentivise additional home owners to get connected need to be identified. For instance, an increase in funding of required residential transmission stations is likely to motivate more owners. However, the funding and therefore expenditure of public means should not be higher than necessary. For this purpose, we propose auto-parameterisation of appropriate measures, such that the intervention target is reached but the funding is minimised.

The task can be interpreted as a specific application of parameter estimation. There have been many techniques applied for agent-based models to determine parameter values that minimize the gap between model results and real world observations. However, the auto-parametrisation is non-trivial because of non-linearity through e.g. agent interaction and micro–macro-level feedbacks. For example, a rapid uptake of heat pump installations may have negative effects on the final diffusion if initial adopters make bad experience (e.g. because of immature technology) and influence their peers in the way they advise them against this technology.

For these reasons, metaheuristics such as genetic algorithms have been applied [15]: Initially, a population of parameter sets (chromosomes) is defined. Their fitness values, i.e. the deviation of simulation results by these parameters and the observation or target state, are calculated. The best scoring chromosomes become parents and are subject to crossover and mutations of genes to produce their offspring. Offspring and parents then form the population of chromosomes for the next iteration until a stopping criterion is reached (see Fig. 5). Still, these approaches require many model runs to compute the fitness function of genetic variations, especially because of many stochastic runs per variation. These are necessary to reflect uncertainty in processes and remaining parameters.

Broniec [16] propose the use of random variable distributions for agent properties and polynomial functions as proxies for their behaviour to substitute ABM

**Fig. 5** Process of applying a genetic algorithm to parameter estimation

Population

Calculation of fitness and Selection

Replacement

Parents

Mutation and Crossover

Offspring

runs during the calculation of fitness scores. Whereas distributions may reflect the heterogeneity of individuals, polynomials determine how their shape and magnitude changes each time step. For basic ecological models such as predator and prey or in epidemiology (SIR model) the method works well. However, for social simulations like investment decisions path dependencies of single agents are expected to be influential. For the application of investments in heating technology, a modified approach is to be developed.

## 3 Discussion

The concise initial literature review of applications of ABM on the adoption of energy efficient technology highlights its potential to foster sustainable transitions of heating technology in the building sector. With the proposed ABM approach, we aim to implement three innovative aspects:

First, we explicitly represent the dynamics evolving from interactions with a broader set of stakeholders, namely energy consultants whose knowledge and thus recommendations shape the set of considered options, and craftspeople such as plumbers whose experience has an impact on building owners' decisions.

Second, we integrate the agent-based model with a model of the local energy system, therefore allowing the representation of feedbacks between the heating infrastructure decisions on the level of municipalities and investment decisions of building owners.

Third, we conduct empirical studies among building owners as well as interviews with municipality representatives, energy consultants and craftspeople. The auto-parameterisation process enables the identification of interventions that allow pre-defined states of residential heat provision. The modelling effort will be accompanied by a series of so-called decision theatres, enabling the feedback of stakeholders about the modelling approach [17]. Results of ABM simulations will be played back to them to validate the model by expert knowledge but also feed the discussion with valuable information about complex interactions.

In existing modelling studies on energy-related adoption decisions, these three innovations have been identified as missing but of crucial importance [4, 6]. Future

work will address the conception and implementation of the multi-stage decision making process involving intermediate phases of information gathering and social influence between trigger and final decision. Interventions which impact particular phases need to be operationalised. Various data sources will be fit to the district level and aligned as a common ground for the ABM and the energy system model.

# References

1. Ehlerding, S.: Mit der Wärmewende aus der Krise, Tagesspiegel Background, https://background.tagesspiegel.de/energie-klima/mit-der-waermewende-aus-der-krise (2022). Last accessed 7/6/2022
2. Senkpiel, C., Dobbins, A., Kockel, C., Steinbach, J., Fahl, U., Wille, F., Globisch, J., Wassermann, S., Droste-Franke, B., Hauser, W., Hofer, C., Nolting, L., Bernath, C.: Integrating methods and empirical findings from social and behavioural sciences into energy system models—motivation and possible approaches. Energies **13**, 4951 (2020). https://doi.org/10.3390/en13184951
3. Ipsos. Beweggründe und Hindernisse für energetische Sanierung (2019)
4. Hesselink, L.X., Chappin, E.J.: Adoption of energy efficient technologies by households—barriers, policies and agent-based modelling studies. Renew. Sustain. Energy Rev. **99**, 29–41 (2019). https://doi.org/10.1016/j.rser.2018.09.031
5. Nava-Guerrero, G.-D.-C., Hansen, H.H., Korevaar, G., Lukszo, Z.: An agent-based exploration of the effect of multi-criteria decisions on complex socio-technical heat transitions. Appl. Energy **306**, 118118 (2022). https://doi.org/10.1016/j.apenergy.2021.118118
6. Du, H., Han, Q., de Vries, B.: Modelling energy-efficient renovation adoption and diffusion process for households: a review and a way forward. Sustain. Cities Soc. **77**, 103560 (2022). https://doi.org/10.1016/j.scs.2021.103560
7. Studer, S., Rieder, S.: What can policy-makers do to increase the effectiveness of building renovation subsidies? Climate **7**, 28 (2019). https://doi.org/10.3390/cli7020028
8. Statistische Ämter des Bundes und der Länder: Gebäude- und Wohnungsbestand in Deutschland: erste Ergebnisse der Gebäude- und Wohnungszählung 2011. Hannover. https://www.statistischebibliothek.de/mir/receive/DEMonografie_mods_00004577 (2014). Last accessed 9/8/2022
9. Arning, K., Dütschke, E., Globisch, J., Zaunbrecher, B.: The challenge of improving energy efficiency in the building sector: taking an in-depth look at decision-making on investments in energy-efficient refurbishments. Energy Behav. **2020**, 129–151. https://doi.org/10.1016/B978-0-12-818567-4.00002-8
10. Zaunbrecher, B.S., Arning, K., Halbey, J., Ziefle, M.: Intermediaries as gatekeepers and their role in retrofit decisions of house owners. Energy Res. Soc. Sci. **74**(101939), 1–12 (2021). https://doi.org/10.1016/j.erss.2021.101939
11. Briegel, R., Ernst, A., Holzhauer, S., Klemm, D., Krebs, F., Martinez Pinánez, A.: Social-ecological modelling with LARA: a psychologically well-founded lightweight agent architecture. In: Seppelt, R., Voinov, A.A., Lange, S., Bankamp, D. (eds.) International Congress on Environmental Modelling and Software 2012, Leipzig, Germany (2012)
12. Ajzen, I.: The theory of planned behavior. Organ. Behav. Hum. Decis. Process. **50**, 179–211 (1991). https://doi.org/10.1016/0749-5978(91)90020-T
13. Sopha, B.M., Klöckner, C.A., Hertwich, E.G.: Adoption and diffusion of heating systems in Norway: coupling agent-based modeling with empirical research. Environ. Innov. Soc. Trans. **8**, 42–61 (2013). https://doi.org/10.1016/j.eist.2013.06.001
14. Krebs, F.: Heterogeneity in individual adaptation action: modelling the provision of a climate adaptation public good in an empirically grounded synthetic population. J. Environ. Psychol. **52**, 119–135 (2017). https://doi.org/10.1016/j.jenvp.2016.03.006

15. Calvez, B., Hutzler, G.: Automatic tuning of agent-based models using genetic algorithms. Lecture notes in computer science, Springer Berlin Heidelberg, 2006, 41–57. https://doi.org/10.1007/11734680_4
16. Broniec, W.: Guiding parameter estimation of agent-based modeling through knowledge-based function approximation. In: Proceedings of the AAAI 2021 Spring Symposium on Combining Machine Learning and Knowledge Engineering (AAAI-MAKE 2021), 2021
17. Wolf, S., Fürst, S., Geiges, A., Laublichler, M., Mielke, J., Steudle, G.: The Decision Theatre Triangle for societal challenges, Global Climate Forum (2021)

# How Beliefs on Food and Climate Change Impact the Dietary Adoption? An Agent-Based Approach

**Maël Franceschetti, Cédric Herpson, and Jean-Daniel Kant**

**Abstract** This paper introduces G-Impact, an agent-based model that combines modelling of household consumption and belief diffusion. Household decisions integrate personal impacts (quality, cost), perceived consequences (climate change, human responsibility), and social norms. The evaluation of these different criteria relies on household beliefs, which can be exchanged during social interactions. These beliefs can be used to explain household decisions on a macro and micro scale, and thus to target information or incentive policies. The model is applied to dietary choice in France, among the omnivorous diet (INCA3), the flexitarian diet and the vegetarian diet. In the control simulation, we observe a significant increase in the proportion of flexitarians, and a slight increase in the proportion of vegetarians over 5 years. We also illustrate the need to properly inform households with the emergence of fake news.

**Keywords** Agent-based model · Consumer behavior · Opinion dynamics · Social norms · Dietary adoption

## 1 Introduction

Given the growing climate risks on a global scale and the involvement of human activities in global warming, it is urgent to act. It is necessary to know which individual and collective behaviors are the most damaging—or virtuous—for the environment, in order to adapt our daily actions. Most climate simulations model human activity in an aggregate way, in the form of a typical behavior representing all individuals.

M. Franceschetti (✉) · C. Herpson · J.-D. Kant
CNRS, Sorbonne Université, LIP6, F-75005 Paris, France
e-mail: mael.franceschetti@lip6.fr

C. Herpson
e-mail: cedric.herpson@lip6.fr

J.-D. Kant
e-mail: jean-daniel.kant@lip6.fr

These approaches do not model the complexity and variability of human behavior or the social interactions that impact individuals' decisions.

To move in that direction, we designed G-Impact, an agent-based model of household consumption, which integrates belief diffusion. Beliefs are the support of product (or service) adoption, and are diffused on a social network, where agents will exchange information on the productions, including their impacts on climate, i.e. greenhouse gas (GHG) emissions. In this paper, we apply G-Impact on dietary adoption, and show how beliefs help us to understand why a particular diet is adopted, and also the impact of fake news. The model is based on real data, found in various surveys among the French population.

## 2    Background

Numerous models are available to simulate the behavior of a population at the macroscopic scale. These models are mainly based on the resolution of differential equations until an equilibrium is reached. The most recent dynamical systems (e.g. Enroads [8]) integrate numerous human activities and their impact on the climate. These approaches test a given scenario with predefined behaviors and do not allow for population heterogeneity or social influences, which nevertheless play a crucial role. These models are not explicable at the individual level.

The agent-based approach allows us to take this heterogeneity into account and to model individuals or households, their decisions and their social interactions. In the BENCH model [6] focused on household energy decision making, households exchange information about their electricity use and make decisions based on personal norms, global warming awareness, and feelings of guilt and responsibility.

The Consumat model [4] is a generic model based on the notions of need and satisfaction. It includes existential and social needs, and personal tastes. An agent who is dissatisfied with his situation will look for alternatives, by informing himself and/or by imitating other agents. This model has been applied to the diffusion of electric vehicles in the STECCAR model [5]. However, the behaviors of the consumat agents may be hard to explain, due to its relative complexity.

An agent-based model has been proposed for meat consumption [7]. It is specific to this topic, do not include the pleasure of eating, nor the flexitarian diet: the model is limited to binary choices (eat or not eat meat).

The model [9] goes further into explicitness, detailing the arguments that drive agents' opinions. Each agent has a set of arguments, represented as an argumentation graph, that they exchange during social interactions. Nevertheless, this approach requires to build a database of arguments as well as the graph of attacks between arguments. It also imposes a particular form of logical reasoning, based on arguments.

Our aim is to propose a flexible model, generic (i.e. to specific to food adoption), easily extensible, and applicable on existing data (e.g. opinion surveys). The resulting model, G-Impact, is inspired by the Theory of Planned Behavior [1] and is based

on household beliefs and their diffusion. Moreover, it includes the notions of global warming awareness, guilt and responsibility.

## 3 G-IMPACT Model

The population is divided into households. The latter have activities to perform and must choose how to perform them. For each activity, they have several options from which they choose: we call it "modalities".

Households have both beliefs and cost estimates about the modalities, allowing them to choose the most appropriate one for each activity. They can be enriched through social interactions, during which households exchange beliefs and information. Households are also sensitive to social norms.

### 3.1 Household Representation

**Household agents** Household agents are representative of the households of the population. Each household has the following attributes:

- the list of household members, along with their attributes: age, gender, employment status, socio-professional category, income,
- the beliefs of the household as well as the estimated costs of the modalities,
- the share of the budget allocated to each type of activity,
- the household's social network,
- an information base, derived from social interactions, indicating the modality choices of other households as well as received beliefs.

**Activities** Households perform activities in different areas, for example: food, transportation, housing, other consumption. These activities are representative of the population's consumption behavior. Each activity is repeated regularly, some daily, others less frequently. Households choose how they will carry out each activity: what modality they will use. Each modality has a (financial) cost, which households estimate, and a real GHG impact, unknown to households.

**Beliefs** A belief is an elementary piece of information, considered to be true. We denote $B_a(o, X)$ the value of a belief for household $a$ between the social object $o$ (product, modality …) and a evaluation criterion $X$. This value is a continuous value in $[-1, 1]$, where 1 denotes a full agreement that $X$ is true for $o$, and $-1$ a full disagreement. For instance $B_a(organic, heathy) = 1$ means that agent $a$ strongly believes that organic food is healthy.

However, each belief has a *reliability* level, to modulate its impact. The reliability depends on its source, from most to less believable: direct experience is most reliable, then comes indirect experience, plausible and then advertising.

Each household has a set of beliefs that allow it to evaluate the modalities and estimate their consequences. We represent beliefs about:

- the modalities: pleasure, impact on health, time saving and GHG reduction[1],
- the climate: the human responsibility in climate change and the perceived impact of global warming on humans.

The beliefs of the households will be represented using associative networks as in the CoBAN [10] model. Each household will have an IAN (Individual Associative Network) representing all its beliefs.

**Estimated costs** Households are provided with a table, containing for each modality the expected cost per individual. A level of reliability is associated with these estimates.

## 3.2 Social Interactions

At each tick, each pair of connected households in the social network have a probability $p_{inter}$ to have a social interaction. These interactions consist of the exchange of messages that may contain beliefs and cost estimates about certain social objects, and the modalities chosen by the household.

Households discuss the climate but also the modalities that interest them out of proactivity (preferred modalities) or out of curiosity (missing information). Each topic of interest has a probability $p_{diffusion}$ to be discussed.

**Integration of beliefs and estimated costs received** When a household receives a message in a social interaction, the beliefs as well as the estimated costs it contains are filtered by decreasing their reliability level.

Then, a belief is directly added to the IAN with a probability $p^{new}$ if it did not exist. Otherwise, there are two cases:

- if the received belief is the same than existing belief and has a higher reliability, then belief is preserved and the maximum reliability level is kept,
- if the received belief contradicts the existing belief, then the existing belief will be revised with some probability, depending on the reliability level of the current and received beliefs (using a probability table).

The same procedure applies to the estimated costs table. The probabilities of revision according to the reliability of the existing belief and that of the new belief are presented in a table, provided as a model parameter. Then, the revision of beliefs is inspired by Assimilation-Contrast theory [2, 3] :

$$c \leftarrow c + sign(r - c) \times \gamma \tag{1}$$

---

[1] We limit ourselves to these four criteria for the moment, which are applicable to all types of consumption and for which we can generally find data.

with $c$ the household's current belief value, $r$ the received belief value, and $\gamma$ drawn uniformly $\in [0, 0.5 - 0.5 \times (|(r - c) - 1|)^{1.3}]$.

Households also aggregate received beliefs into a normative IAN, used to construct the average perceived social representation of each modality.

## 3.3 Household Cognition

**Evaluation of beliefs** Households can evaluate their beliefs. This mechanism allows them to extract from the IAN the value of a belief between a social object and a desired concept. A belief must have a reliability higher than $\sigma^{threshold}$ to be evaluated: information that is too unreliable is not taken into account.

**Update of estimated costs** If a household does not have an estimate of the cost of the modality it is evaluating, it will fill in this missing information using the mean cost estimate of the other modalities for this activity.

When a household adopts a modality, its estimated cost is updated:

- when no estimate of the cost is available, it is initialized as a random number around the real price, with a maximum percentage of error $max\_init\_error$,
- when an estimate is available, it is adjusted to be closer to the real price (using a linear interpolation between the current estimate and the real cost, where the adjustment coefficient is drawn uniformly in [0,1]).

**Decision making** We define the following notations:

- $a$ the household concerned (the one which evaluates the modality),
- $act^m$ the modality $m$ of the activity $act$,
- abbreviations: $CC$ for Climate Change, $HR$ for Human Responsibility, $PCC$ for Perceived Climate change Consequences, $GR$ for GHG Reduction.

Households assess the usefulness ($U$) of each modality according to the criteria of personal impacts ($PI$), the perceived consequences ($PC$), and social norms ($N$). The chosen modality is the one with the highest utility value. Changing modality is binding, a resistance to change factor is applied in this case.

Personal Impacts ($PI$) are related to costs and quality. A modality must be within budget and prices affect decisions. The quality of a modality is assessed based on beliefs about the pleasure, health, and time-saving.

Perceived consequences ($PC$) are related to climate change consequences and household responsibility. The more a household feels responsible in climate change, the less it will positively evaluate a modality associated with a weak reduction in GHG emissions.

Norms ($N$) considered are descriptive and injunctive norms: the more a modality is used and appreciated by the population, the more it is considered socially accepted.

$$U_a(act^m) = ca(act^m) \times \left(1 + \frac{PC_a(act^m) + PI_a(act^m) + N_a(act^m)}{3}\right) - 1 \quad (2)$$

with $ca(act^m) \in [0, 1]$ the change acceptance factor, $ca(act^m) = 1$ if $act^m$ is currently used, else $ca(act^m) = (1 - \rho)$; $\rho \in [0, 1]$ the resistance to change factor.

**Perceived consequences *(PC)***

$$PC_a(act^m) = B_a(act^m, GR) \times AR_a(act^m) \times B_a(CC, PCC) \quad (3)$$

*Estimated responsibility of a household (AR)*

$$AR_a(act^m) = \left(\frac{B_a(CC, HR) + RR_a(act^m)}{2} + 1\right) \times \frac{1}{2} \quad (4)$$

*Relative responsibility of a household (RR)*

$$RR_a(act^m) = \begin{cases} 0 \text{ if } nb_a = \emptyset \\ \dfrac{\sum_{n \in nb_a}(C_a(act_{mod(n,act)}) - C_a(act^m))}{2 \times |nb_a|} \text{ else} \end{cases} \quad (5)$$

**Personal impact *(PI)***

$$PI_a(act^m) = CU_a^{sp}(act^m) \times QU_a^{sq}(act^m) \quad (6)$$

with $sp \in [0, 1]$ the sensitivity of PI to the price, $sq \in [0, 1]$ the sensitivity of PI to the quality (depends on the household category).

*Cost utility for a household (CU)*

$$CU_a(act^m) = \exp\left(\frac{-price_a(act^m)}{(1 - bp) \times max\_price}\right) \quad (7)$$

with $price_a(act^m) \in R^+$ the net (including any subsidies) estimated costs of $act^m$ for the household; $bp$ the households sensitivity to low prices; $max\_price$ the maximum estimated price among all possible modalities for this activity. If the price is out of budget, we give $CU$ the value $-\infty$.

*Quality utility for a household (QU)*

$$QU_a(act^m) = \frac{1}{|CR|} \times \sum_{i \in CR} \frac{1 + B_a(act^m, i)}{2} \quad (8)$$

with $CR = \{health, \ time \ saving, \ pleasure\}$, the set of criteria for quality.

**Social norms** *(N)*

$$N_a(act^m) = \frac{DN_a(act^m) + IN_a(act^m)}{2} \tag{9}$$

*Descriptive norm perceived by a household for a modality (DN)*

$$DN_a(act^m) = \begin{cases} 0 \text{ if } nb_a = \emptyset \\ \left( \frac{2}{|nb_a|} \times |\{n \in nb_a , \; mod(n, act) = m\}| \right) - 1 \text{ else} \end{cases} \tag{10}$$

with $mod(n, act)$ the modality chosen by household $n$ for action $act$, and $nb_a$ the neighboring households of $a$ in the social network.

*Injunctive norm perceived by a household for a modality (IN)*

$$IN_a(act^m) = \frac{NQU_a(act^m) + NPC_a(act^m)}{2} \tag{11}$$

*Normative quality utility (NQU)*

$$NQU_a(act^m) = \frac{1}{|CR|} \times \sum_{i \in CR} \frac{1 + NB_a(act^m, i)}{2} \tag{12}$$

with $NB_a(o, x)$ the evaluation of the belief between the social object $o$ and the concept $x$ in the normative IAN.

*Normative perceived consequences (NPC)*

$$NPC_a(act^m) = NB_a(act^m, GR) \times \frac{NB_a(CC, HR) + 1}{2} \times NB_a(CC, PCC) \tag{13}$$

## 3.4 Course of the Simulation

The simulation life cycle is decomposed into three steps: initialization, execution of each time step (tick), and ending. At each time step, the course is as follows:

1. Choice of modalities: for each activity to be performed, households evaluate all the modalities and choose which one they will use.
2. Execution of the modalities: all households execute the chosen modalities. Agents state and GHG emissions are collected.
3. Social interactions: households can interact with each other.
4. Beliefs and estimates update: households process all messages received and update their beliefs, cost estimates and normative IAN.

## 4    Application to Dietary Adoption

### 4.1    Actions and Modalities

We applied the G-Impact model to food consumption in France. Households must choose what they eat, i.e. their diet: INCA3 (most common diet in France, omnivorous with meat), flexitarian (meat reduction) or vegetarian (no meat). This decision is made once a week and the decision applies for the entire week.

### 4.2    Initialization of the Population and Beliefs

The strength of our instantiation of the model is to use as much data as possible on the French population, mainly between 2017 and 2019.

The cost and GHG emissions of the modalities were taken from different national and international studies (ANSES, INSEE, WWF).

Population and income are initialized from national INSEE data. Households are assigned price and quality sensitivity values ($sp$, $sq$, $bp$) according to their income and the individuals that compose them. Our population of 9943 household agents is representative of the French population. To initialize the beliefs of the households, we convert Likert scales taken from several national opinion surveys into values of beliefs in $[-1, 1]$, which we distribute according to the proportions indicated. This information is supplemented with national studies when necessary. In the data used for the experiments, the initial beliefs of households about their own diet have the same distributions regardless of the diet, only the assumptions about diets other than theirs vary. We do not give the INCA3 dieters any preconceived notions about the impact of the vegetarian and flexitarian diets on health and the environment: we want to study how, during the simulation, vegetarians and flexitarians manage to spread these beliefs to the whole population.

We draw for each household its initial diet according to the declared distribution in the real population. We add a Gaussian noise on the initial beliefs of the households.

Once the households have been generated, we create the social network of the population, using a Small-World [11] network linking the households together.

### 4.3    Experiments

**Control simulation** We ran 30 simulations with a basic set of parameters over 5 years. We used $p_{inter} = 0.02$, $p_{diffusion} = 0.05$, $\rho = 0.08$, $p_{new} = 0.9$, $max\_init\_error = 20\%$, $\sigma^{threshold} =$ 'advertising'. The filtering of received beliefs reliability is the following: direct experience become indirect experience, and indirect experience become plausible. With respect to the reliability of the current belief, the probabili-

**Table 1** Distribution of the different diets at the start and at the end of the simulation (mean and std)

| Diet | Initial | Final |
|---|---|---|
| INCA3 | 74.46% (± 0.45) | **63.41**% (± 1.39) |
| Flexitarian | 20.35% (± 0.50) | **28.59**% (± 1.48) |
| Vegetarian | 5.18% (± 0.20) | 7.99% (± 0.63) |

**Table 2** Proportion of final diet for each initial diet (mean and std)

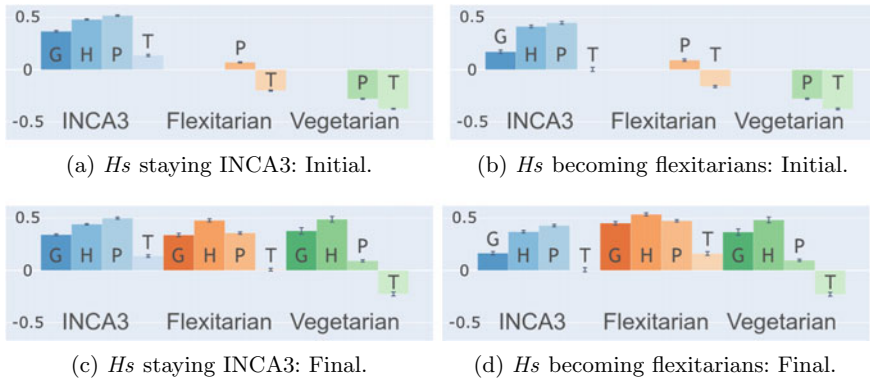| | | Final diet | | | |
|---|---|---|---|---|---|
| | | INCA3 | Flexitarian | Vegetarian | Total (%) |
| Initial diet | INCA3 | **73.47**% (± 1.24) | 20.41% (± 1.28) | 6.12% (± 0.53) | 100 |
| | Flexitarian | 33.56% (± 1.47) | **59.89**% (± 1.77) | 6.55% (± 0.79) | 100 |
| | Vegetarian | 36.15% (± 2.98) | 23.22% (± 2.57) | **40.64**% (± 2.75) | 100 |

**Table 3** Households average beliefs depending on final diet (mean and std)

| | GHG reduction | Health | Pleasure | Time saving |
|---|---|---|---|---|
| INCA3 | 0.317 (± 0.008) | **0.387** (± 0.008) | **0.457** (± 0.011) | 0.133 (± 0.013) |
| Flexitarian | **0.436** (± 0.015) | **0.527** (± 0.012) | **0.486** (± 0.011) | 0.165 (± 0.015) |
| Végétarian | **0.473** (± 0.030) | **0.541** (± 0.025 ) | 0.261 (± 0.024) | 0.043 (± 0.031) |

ties of revision are 1.0 if the received reliability is higher, 0.9 if it is equal, 0.01 if it is just below, 0.001 if it is even lower. Households have an average of 10 neighbors in the social network.

We see in Table 1 that the proportion of practitioners of the INCA3 diet has significantly decreased at the end of the simulation ($-11.05$% points), in favor of flexitarian and vegetarian diets ($+8.24$ and $+2.81$% points respectively), leading to a reduction of annual GHG emissions from food associated with diet choice of 5% between the first and the last year. We can see in Table 2 that a significant proportion of vegetarians and flexitarians eventually became practitioners of the INCA3 diet (respectively 33.56% and 36.15%). The proportion of vegetarians who have maintained their diet is only 40.64%, thus representing the difficulty in maintaining this diet.

*What are the belief profiles of adopters of different diets?* We display the average final beliefs of households about the diet they follow in Table 3. Following the diffusion of household beliefs, three belief profiles emerge corresponding to the three diets. There are two major criteria for INCA3 diet, which are pleasure and health. Flexitarian diet has three more homogeneous criteria, which are health, GHG reduction and pleasure. Vegetarian diet has two major criteria: health and GHG reduction. The choice of the INCA3 diet is then more associated with "selfish" criteria, flexitarian diet on a more multi-criteria and balanced decision, and for the vegetarian diet it is a more "altruistic" choice.

(a) *Hs* staying INCA3: Initial.



(b) *Hs* becoming flexitarians: Initial.



(c) *Hs* staying INCA3: Final.
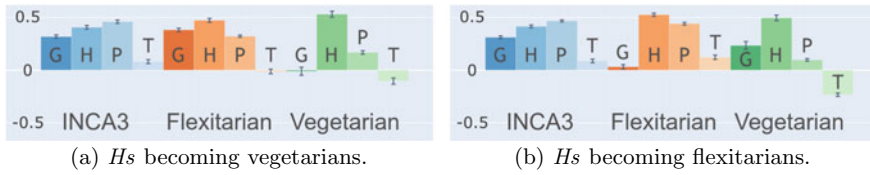


(d) *Hs* becoming flexitarians: Final.

**Fig. 1** Average initial and final beliefs (and std over the simulations) of households (*Hs*) initially INCA3 dieters. For each diet, beliefs are: GHG reduction (G), health (H), pleasure (P) and time saving (T)

*What caused these transitions from the INCA3 diet to the flexitarian diet?* We observe in Fig. 1d that the households initially practicing the INCA3 diet that have finally adopted a flexitarian diet had a rather poor opinion about the INCA3 diet compared to those who maintained this diet, in particular its ability to reduce GHG emissions and the time savings associated with it. Households that finally adopted the flexitarian diet also received beliefs with higher values about it: they have been more convinced of the benefits of the flexitarian diet.

**Diffusion of a fake news** We are now going to illustrate the impact that the beliefs of a minority of the population can have on the whole of the latter, via the dissemination of fake news. We attribute to all households an extreme value of $-1$ for the belief on the perception of the impact of global warming (PCC): households thus believe that the latter will only have positive effects on humans.

We then observe that a minority of vegetarians and flexitarians, who believe that their diet does little to reduce GHG emissions (and therefore that it allows global warming, perceived as favorable here), very strongly disseminate these beliefs out of proactivity. This is illustrated in Fig. 2, where households initially INCA3 practitioners received beliefs close to 0 concerning the reduction of GHG emissions from the vegetarian or flexitarian diet that they finally chose. The impact of this diffusion is significant: at the end of the simulation 34.5% of the population is flexitarian and 19.5% is vegetarian, ironically leading to a reduction of GHG emissions from food associated with diet choice of 14% compared to the control simulation (comparison over the last year). It should be noted that these beliefs spread more easily about minority diets, for which less contradictory beliefs that could stop these fake news circulate. In this experiment, the fake news (vegetarian and flexitarian diets do not help to reduce GHG emissions) entails a virtuous choice, but for a wrong and paradoxical reason, when some essential information is missing (climate change is not good).

(a) *Hs* becoming vegetarians.   (b) *Hs* becoming flexitarians.

**Fig. 2** Average final beliefs (and std over the simulations) of households (*Hs*) initially INCA3 dieters who have adopted another diet. For each diet, beliefs are: GHG reduction (G), health (H), pleasure (P) and time saving (T)

## 5 Discussion

In this paper, we have presented the outline of G-Impact, an agent-based model that combines modelling of household consumption and belief diffusion. The method used to apply the model to food can be followed to apply the model to other types of consumption (simultaneously or not).

The strength of our approach is that it provides a descriptive but also explanatory analysis, notably via beliefs, at the macro and micro levels. This would allow the design and testing of information and incentive policies to reduce GHG emissions, in particular policies targeting specific population groups.

Preliminary experiments highlight the important impact of the diffusion of beliefs and social interactions on household consumption behavior. We have seen that a minority of the population, ill-informed about their own practices, can spread false information to a large part of the population. It is therefore very important to inform the population widely about their own practices, even when they are in the minority.

Several elements of the model can still be improved, such as adding weights on the criteria, incorporating personal norms and ethical criteria, or allowing households to acquire beliefs after the choice of modalities via a feedback loop.

## References

1. Ajzen, I.: The theory of planned behavior. Organ. Behav. Human Decis. Proc. **50**, 179–211 (1991)
2. Freedman, J.L.: Involvement, discrepancy, and change. J. Abnormal Soc. Psychol. **69**(3), 290 (1964)
3. Hovland, C.I., Harvey, O., Sherif, M.: Assimilation and contrast effects in reactions to communication and attitude change. J. Abnormal Soc. Psychol. **55**(2), 244 (1957)
4. Jager, W., Janssen, M.: An updated conceptual framework for integrated modeling of human decision making: the consumat ii. In: Paper for Workshop Complexity in the Real World@ ECCS, pp. 1–18 (2012)
5. Kangur, A., Jager, W., Verbrugge, R., Bockarjova, M.: An agent-based model for diffusion of electric vehicles. J. Environ. Psychol. **52**, 166–182 (2017)
6. Niamir, L., Filatova, T., Voinov, A., Bressers, H.: Transition to low-carbon economy: assessing cumulative impacts of individual behavioral changes. Energy Policy **118**, 325–345 (2018)

7. Scalco, A., Macdiarmid, J.I., Craig, T., Whybrow, S., Horgan, G.W.: An agent-based model to simulate meat consumption behaviour of consumers in Britain. J. Artif. Soc. Soc. Simul. **22**(4), 8 (2019)
8. Siegel, L.S., Homer, J., Fiddaman, T., McCauley, S., Franck, T., Sawin, E., Jones, A.P., Sterman, J., Interactive, C.: En-roads simulator reference guide. Tech. Rep. (2018)
9. Taillandier, P., Salliou, N., Thomopoulos, R.: Introducing the argumentation framework within agent-based models to better simulate agents' cognition in opinion dynamics: application to vegetarian diet diffusion. J. Artif. Soc. Soc. Simul. **24**(2) (2021)
10. Thiriot, S., Kant, J.D.: Using associative networks to represent adopters' beliefs in a multi-agent model of innovation diffusion. Adv. Complex Syst. **11**(2), 261–272 (2008)
11. Watts, D.J., Strogatz, S.H.: Collective dynamics of 'small-world' networks. Nature **393**(6684), 440–442 (1998)

# Public Acceptance of Green Mobility Policies: An Agent-Based Model


Check for updates

**Marie Lisa Kogler, Annina Thaller, and Daniel Reisinger**

**Abstract** We present an agent-based model to simulate policy acceptance for push and pull policy measures. Push measures are generally perceived as restrictive and are often directed towards the reduction of private car use, e.g. fuel price increases and inner-city car bans. Pull measures relate to diverse incentives to facilitate climate-friendly travel choices, e.g. attractive offers for public transport such as interregional cost reductions and expansion of the public transport infrastructure. The model is informed by empirical data regarding agents' travel mode utilities and allows to evaluate agents' satisfaction and acceptance of diverse policy scenarios. Regional dependencies are tested for the case of Austria. The results show that the political acceptance of push measures increases when they are combined in packages of measures considering the expansion of public transport infrastructure. Furthermore, the general acceptance of green mobility measures is closely linked to the existing infrastructure of the individual districts.

**Keywords** Climate mitigation · Agent-based · Transport · Public opinion · Emission reduction

## 1 Introduction

As the latest IPCC report [20] highlights, transport is one of the critical areas in mitigating climate change. Globally, transport currently accounts for 15% of total greenhouse gas (GHG) emissions and 23% of energy-related $CO_2$ emissions. "*Transformative changes*" are needed to meet climate targets. Without such changes, transportation $CO_2$ emissions are projected to increase by 16 to 50% by 2050 [20]. While the challenge is clear, the necessary behavioral changes that would enable low-carbon mobility remain elusive. One reason is that self-driven behavior change in passenger transport is unlikely due to various sociocultural and institutional factors that act,

M. L. Kogler (✉) · A. Thaller · D. Reisinger
Institute of Environmental Systems Sciences, University of Graz, Merangasse 18, 8010 Graz, Austria
e-mail: marie.kogler@uni-graz.at

511

at least in part, as barriers to demand reduction and modal shift [20, 25]. At the same time, an implementation gap [4] in stringent policy measures to drive such needed behavior change is apparent [24]. A common obstacle in this regard is the public acceptance of such measures, or lack thereof, which can act as both an enabler and a barrier to passenger transport transformation [26]. To achieve the necessary emission reduction targets, it is essential to use a variety of policy instruments, often referred to as policy packages [7]. Such policy packages should aim for a combination of restrictive, i.e. push, and incentive-creating, i.e. pull, measures in order to achieve both high and rapid effectiveness and implementability of the proposed measures [23]. Push measures involve various restrictions, such as fuel price increases or driving bans, while pull measures try to make alternatives more attractive, e.g. by improving public transport or bicycle infrastructure [18]. Typically, pull measures are preferred by the public and policymakers over push measures because they are less intrusive and therefore easier to implement without much resistance [5]. At the same time, however, it is precisely such restrictive measures that are essential for successful climate change mitigation [3]. Therefore, it is important to obtain broad public support in order to implement effective policies and policy packages.

Within the context of sustainable passenger transport, agent-based models (ABM) are used to simulate public commuting behavior and modal shifts, often by scenario analysis to test the impact of different policies [1, 2, 8, 11, 22]. Systematic review of ABM in mobility transition [16] shows that there is a strong focus on i) electric mobility [17, 21, 27], and alternative fuel vehicles (e.g., hydrogen cars) [12, 13] including charging and refueling behavior, and ii) urban related mobility options such as micromobility (e.g., electric bicycles and scooters), and mobility services (e.g., car sharing, bike sharing, bus services, shared automated vehicles). An notable contribution regarding urban mobility and simulation of low carbon commuting choices is provided by Maggi and Vallino [15].

While the above ABMs regarding mobility behavior provide insight into estimating environmental impacts as a result of certain behaviors of actors [14], the feasibility of implementing the applied policy measures is often neglected. To address the highly relevant issue of public support in passenger transport regulation, this paper investigates public attitudes toward different push and pull measures and their combination into policy packages using an ABM. In particular, we aim to investigate the question of which policy package promises the greatest public support, focusing on two pull and two push measures that are relevant for the Austrian context. The model is informed by geospatial data concerning the usage behavior of cars and public transport. Data driven parameterization is applied to calibrate the model output regarding the satisfaction and acceptance of various proposed measures.

**Table 1** Overview of the selected push and pull policies

| Abbr. | Type | Policy | Wording |
|---|---|---|---|
| A | Push | Fuel price increase | An increase in the taxation of fossil fuels is to be implemented. As part of this increase, the taxation of gasoline and diesel fuel is to be harmonized in order to abolish the existing "diesel privilege" |
| B | Push | Inner-city car ban | An inner-city driving ban for all private cars in cities above a certain size is to be introduced. The ban will apply to all private cars, including e-cars. There are exceptions, for example for residents and loading activities |
| C | Pull | Public transport cost reduction | Financial support for public transport and tax reductions for public transport tickets |
| D | Pull | Public transport expansion | Public transport that is timed according to the possible number of users and also ensures a minimum service for rural areas |

## 2 Data and Methods

### 2.1 Empirical Background

Empirical data used to model public acceptance of various transport policies stems both from primary and secondary data. With regard to primary data, we use the data sets from two quota-representative surveys conducted for the general public of Austria in May 2021 and November 2021. Both consider with different transport policies and respondents' acceptance of such policy proposals. We selected two push measures that target motorized individual transport, one relevant on a national scale, (A) *fuel price increase*, and one specific for urban areas, (B) *inner-city car ban*. In terms of pull measures, we selected two measures that directly relate to public transport, namely (C) *public transport cost reduction* and (D) *public transport expansion*. Find an overview of the selected policies in Table 1.

In Study 1 ($N = 1{,}032$), information was collected on policy $D$, environmental awareness, and mode choices for car travel and public transport. Study 2 ($N = 1{,}084$) collected information on policies $A$ and $B$. In both cases, information was grouped in residential areas, namely urban, suburban, subrural, and rural, based on self-reported answers.[1] For policy $C$, we used secondary data from the literature based on the German study by Engler et al. [6] and on the status report of Austrian mobility [9], as we did not include this policy measure in the aforementioned studies. In Table 2 is an overview of the variables (means and standard deviations) used as input data for the ABM. Other empirical information on policy acceptance was used for evaluation of simulation results.

---

[1] Original answer categories of study 1 and 2 had to be adjusted to be comparable.

**Table 2**  Means $\pm$ standard deviations from the empirical data grouped by residential area

|                        | Urban          | Suburban       | Subrural       | Rural          |
|------------------------|----------------|----------------|----------------|----------------|
| Environment awareness  | $3.51 \pm 0.879$ | $3.36 \pm 0.791$ | $3.40 \pm 0.788$ | $3.41 \pm 0.824$ |
| Car                    | $2.46 \pm 1.19$  | $3.18 \pm 1.12$  | $3.37 \pm 0.898$ | $3.38 \pm 0.897$ |
| Public transport       | $2.50 \pm 1.29$  | $1.64 \pm 0.912$ | $1.66 \pm 0.974$ | $1.40 \pm 0.719$ |

Mode choice for car and public transport ranges on scales from 1 (never) to 5 (daily). Environmental awareness ranges on a scale from 1 (very low) to 5 (very high)

## 2.2  Model

The satisfaction function $F$ for each agent $j$ is defined by

$$F^j = \alpha F^j_{soc} + (1 - \alpha) F^j_{pers} \tag{1}$$

where $F^j \in [0, 1]$. A random variable $\alpha \in [0, 1]$ serves as weight between the social and personal needs. The social needs $F_{soc}$ encompass influences from the individual neighborhood of each agent, given by the adjacency matrix $A$ of the graph. The social influence is represented by the average satisfaction of the agent's neighborhood:

$$F^j_{soc} = \frac{\sum_i^{k_j} F^i}{k} \tag{2}$$

with nearest neighbor $i$ of an agent and the degree of a node $k_j = \sum_j a_{ij}$ of the adjacency matrix $A$, representing the number of nearest neighbors. Note that $F^j_{soc} \in [0, 1]$ and evolves over time.

The personal needs are fueled by empirical insights, following the basic assumption that the mode choice represent the respective utility of a mode of transportation for an agent. Thus, the personal needs consist of a sum of the utility of the different travel modes, in the here presented case on the utility for car $u^j_c$ and utility for public transport $u^j_p$:

$$F^j_{pers} = \frac{u^j_c + u^j_p}{2} \tag{3}$$

with $F^j_{pers} \in [0, 1]$. Empirical studies, as referred to in Sect. 2.1, have shown a linear correlation with a person's environmental awareness and mobility policy acceptance. The higher the environmental awareness, the greater the willingness to support a policy, even if one's own satisfaction is not well met. Thus, the agent's acceptance $H^j$ of a policy or policy package is a binary choice given by:

$$H^j = \begin{cases} 1 & \text{, if } F^j \geq T \\ 0 & \text{, if } F^j < T \end{cases} \tag{4}$$

with the threshold $T = 1 - env$. Through this approach, it follows that environmental awareness acts like a buffer, since $F^j + env \geq 1$. For agents with a high satisfaction, a low environmental friendliness is sufficient to achieve acceptance to a policy. Conversely, if the agent's satisfaction is low, but the environmental awareness is high, acceptance can also be achieved.

**Direct and indirect policy impact** The considered policy measures $m = A, B, C, D$ have different impacts on the utilities of the travel modes. Policy $A$ (fuel price increase) and policy $B$ (inner-city car ban) have a *direct impact* on the utility for car usage $u_c$. The penalty of price increases is proportional to the amount of car usage, while the penalty of inner-city car bans is mainly for short distance travel in cities. Policy $C$ (public transport cost reduction) and $D$ (public transport expansion) have a direct impact on the utility for public transport $u_p$. The benefit of cost reductions is proportional to the typical distance traveled, while the purpose of expanding public transport infrastructure is to facilitate usage for people that have so far low usage behavior.

Next to direct impacts of push and pull policies, there can also be *indirect* effects on other modal choices, which are in competition with each other. In that perspective, push policies $A$ and $B$ have a direct effect on the car utility and an indirect effect on the public transport utility, which becomes more attractive as travel choice. Conversely, pull policies $C$ and $D$ have a positive effect on public transport utility while making car use less attractive. These synergies are particularly important for the design of policy packages and implemented by using an approach based on the Lotka-Volterra model. These relationships were formally implemented as utility updates given as follows:

$$u_c = u_c + u_c^m - \gamma_{cp}^m u_p \tag{5}$$
$$u_p = u_p + u_p^m + \gamma_{pc}^m u_c \tag{6}$$

To simplify notation, the agent index $j$ is omitted from the equations. Direct effects of single policy measures $m$ are given as follows:

$$u_c^A = -\beta^A u_c \tag{7}$$
$$u_c^B = -\beta^B 1/u_c \tag{8}$$
$$u_p^C = \beta^C u_p \tag{9}$$
$$u_p^D = \beta^D 1/u_p. \tag{10}$$

The parameters $\beta$ denote direct policy impacts, while the parameters $\gamma$ denote indirect policy impacts.

**Table 3** Model parameters: Expected value of the social-personal weight $\alpha$, direct $\beta$ and indirect $\gamma$ impact parameters

| Model parameter | A | B | C | D |
|---|---|---|---|---|
| $E(\alpha)$ | 0.5 | 0.5 | 0.5 | 0.5 |
| $\beta$ | 0.4 | 0.05 | 0.1 | 0.6 |
| $\gamma$ | 0 | 0 | 0 | 0 |

**Table 4** Weights of single policy measures in different policy packages for all considered policy scenarios

| Policy packages | A | B | C | D |
|---|---|---|---|---|
| A + C | 0.75 | | 0.25 | |
| A + D | 0.5 | | | 0.5 |
| B + D | | 0.25 | | 0.75 |
| A + B + C + D | 0.3 | 0.2 | 0.2 | 0.3 |

**Modeling policy packages** Policy packages are designed to combine restrictive measures with incentives. Studies show that fuel price increases and public transport expansion in policy packages have a greater effect on perceptions of the package as a whole than other measures. Therefore, the personal needs $F_{pers}$ for policy packages are expanded to a weighted sum of single policy impacts. The weighting of policies in the policy packages $A + C$, $A + D$, $B + D$ and $A + B + C + D$ is shown in Table 4 according to expert opinion (Table 3).

**Regional representation** Additional to the algorithm presented in Sect. 2.2, each agent is associated with a district (116 districts, including 23 districts of the capital Vienna), and a residential area. The classification in residential area allows to regionalize the input data, see Table 2. The classification bounds for different districts are chosen for the case of the Austrian situation: urban above $100k$, suburban $10k - 100k$, subrural $2k - 10k$, rural below $2k$. Geographic data on population distributions such that agents' regional attributes were designed to be representative for the Austrian case.

## 3 Setup and Simulations

The setup of the ABM is as follows: A networked agent population ($N = 841$), which is based on the Austrian population distribution in terms of the assignment of residential areas and districts. The underlying network topology is essentially the Barabási-Albert (BA) graph, a standard network for social interactions, which follows a power-law degree distribution and is therefore scale-free and a good choice for a representative population for a country. Initial variables relevant for the model dynamics
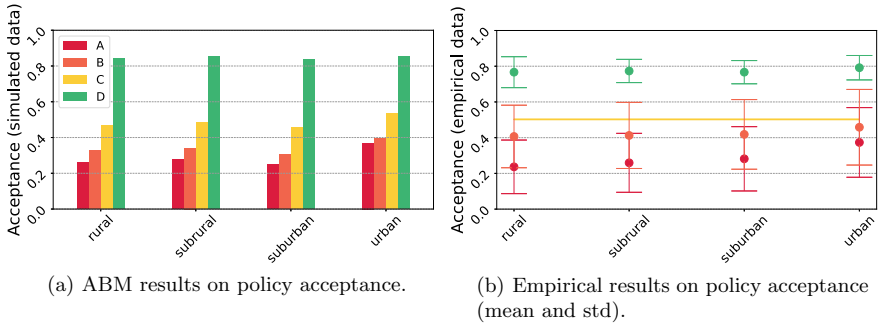
($u_c$, $u_p$ and $env$) are set based on Table 2. The initial value of the agents' satisfaction $F$ represents the current situation of each agent (status-quo). The simulation starts with an assessment of the agents on the respective policy scenario. Subsequently, agents update their satisfaction and acceptance due to social interactions. All simulations end when equilibrium is reached. The scenarios considered refer to single policies $A$, $B$, $C$, $D$ and policy packages $A + C$, $A + D$, $B + D$, $A + B + C + D$. Ensemble simulations of 100 per policy scenario were performed for the model analysis. The agents' associated residential areas and districts are used for later geospatial evaluation of the model output (comparing empirical and simulated data on policy acceptance).

In this first investigation of the presented ABM, emphasis was placed on recognizing and assessing the basic properties of the model dynamics, mediated by the direct policy impact given by $\beta$, see Table 3. The ordinal scaling (ranking) of these values was selected by expert opinion, and a finer adjustment was subsequently made by parameter fitting using the model. To show the general feasibility of the model approach, other influences are not considered within this work. Therefore, the weighting between social and personal needs $\alpha$ is uniformly distributed with an expected value $E(\alpha) = 0.5$, and indirect policy effects were omitted ($\gamma = 0$). More details to this limiting choices are discussed in Sect. 5.

Throughout the model development, simulation phase and analysis, the computational model was thoroughly validated [10, 19]. The validation procedure includes graphical validity (display of micro and macro dynamics), time-line tracing, extensive stochastic comparison to track inconsistencies (internal validity), and sensitivity analysis. The sensitivity of the model is tested via parameter variation. The model output lies in the insights of agents' satisfaction and acceptance in regard to regional differences. The model validation phase also includes testing the results for different network types and variation of network parameter. Ensemble simulations were performed on scale free topologies, small-world topologies and random networks.

## 4 Results

Overall, the model performance showed great alignment with social studies and expert opinions on the topic of policy acceptance for green mobility. This is particularly noteworthy since the model is driven by empirical evidence and only requires few additional assumptions. An important aspect was the agent-based design, as this allowed to model social structures by using suitable interaction topologies. Scale-free and small-world networks performed very well, and results were robust to variations in node degrees. For both topology types, there was good agreement with empirical data, which can be explained by the respective network properties such as strong clustering and other common traits of social structures. Random networks showed an overall low mean acceptance of all policy measures, including pull measures, which are known to be of general high acceptance. Therefore, random networks were not a

(a) ABM results on policy acceptance.

(b) Empirical results on policy acceptance (mean and std).

**Fig. 1** Average acceptance of push ($A$, $B$) and pull ($C$, $D$) policy measures in different residential areas
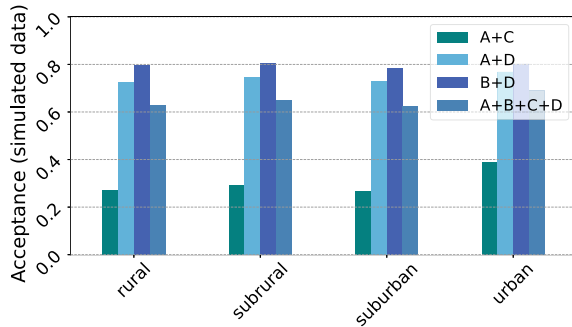
suitable choice for this model type. This is to be expected, since random topologies do not accurately reflect social diffusion processes.

## 4.1 Acceptance of Single Policies and Policy Packages

The comparison of the model results with empirical study data is shown in Fig. 1. For all four of the single policy measures $A$, $B$, $C$, $D$, excellent overlap is seen between the two data sets (mean squared error: 0.0286) and all model results are within the range of confidence (error bars in Fig. 1b) of the available study data. Disparities between the residential areas are well represented by the ABM: the higher acceptance of both push measures in urban areas is reflected by simulation results. Regional variations in the overall high level of acceptance of the expansion of public transport are also clearly visible. Moderate deviations from the ABM data and empirical data are also apparent: policy $B$ shows an overall underestimation of acceptability (mean squared error: 0.05581). Note that in the empirical data of policy $C$, the resolution to residential areas was not available, and the national mean is indicated by a line in Fig. 1b.

Figure 2 gives an overview of the influence on acceptance when policies are combined into packages. Altogether, the least preferred package representing a sole pricing approach is $A + C$ (overall mean: $0.29 \pm 0.041$). The most preferred package is $B + D$ (overall mean: $0.8 \pm 0.023$, combining one push and one pull measure. The combination of all four measures $A + B + C + D$ (overall mean: $0.639 \pm 0.03$) is the only scenario examined in which two push measures are included. This results in a significantly dampened acceptance compared to the scenarios $A + D$ (overall mean: $0.738 \pm 0.025$) and $B + D$. Except for $A + C$, all modeled policy packages attain an acceptance level of over 0.6. This indicates that policy $D$ is a good lever to get broad popular support for a mobility package.

**Fig. 2** ABM results on average acceptance of policy packages $A + C$, $A + D$, $B + D$ and $A + B + C + D$
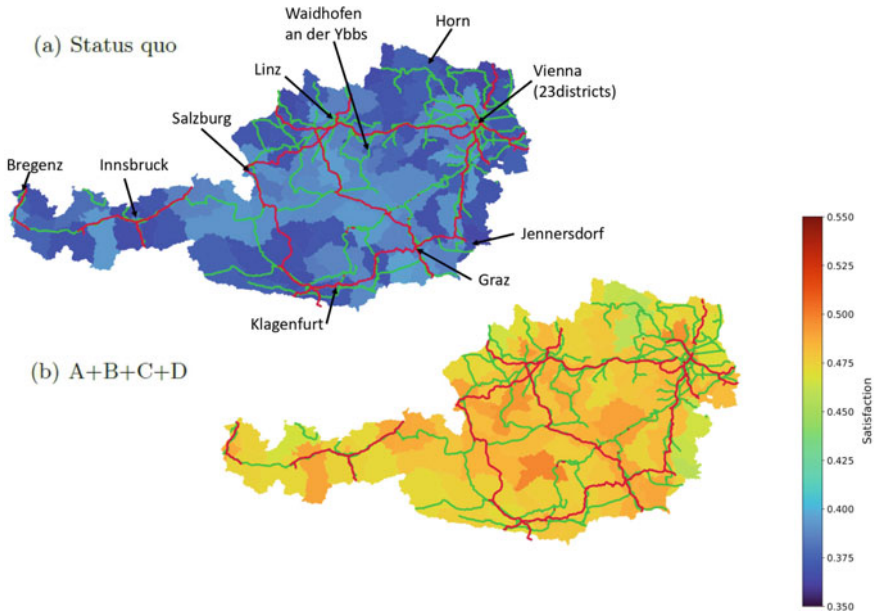
## 4.2 Differences in Satisfaction and Acceptance

For single policies, simulated satisfaction ranges between $0.271 \pm 0.011$ for policy $A$ and $0.655 \pm 0.01$ for policy $D$, representing a decrease in average initial satisfaction for policy $A$ and $B$ (push measures). The acceptance of single policies ranges between $0.275 \pm 0.04$ for policy $A$ and $0.845 \pm 0.019$ for policy $D$. Due to the nature of the mathematical equations used in the model, acceptance is always equal or higher to the satisfaction level. The bigger the differences, the stronger the effect of environmental awareness. This difference is highest for the policy combination $B + D$ (satisfaction = 0.607, acceptance = 0.8), and in general very high as soon as policy $D$ is included.

## 4.3 Regional Results for Districts

In terms of policy $A$, the lowest 10 % of satisfaction can be found exclusively in rural (the lowest to be found in Burgenland, district *Jennersdorf*), and the highest 10 % exclusively in urban districts (the highest in Vienna, district *Liesing*). Regarding policy $D$, the lowest 10 % of satisfaction can be found in suburban and urban areas (minimum in Lower Austria, district *Waidhofen an der Ybbs*). The highest 10 % are mostly located in subrural areas (maximum satisfaction in rural Carinthia, district *Feldkirchen*).

Figure 3 compares the satisfaction without any policy (status quo, left-hand side) with the impact of the policy package $A + B + C + D$ (right-hand side). Overall, satisfaction increases in all districts, on average from $0.378 \pm 0.009$ to $0.48 \pm 0.008$. High values in satisfaction for the full policy package scenario cannot be linked to one specific type of residential area, as they are found in different types as well as numerous federal states. The highest policy satisfaction can be found in inner-city Vienna, the lowest in rural Lower Austria (district *Horn*).

**Fig. 3** Satisfaction **a** without and **b** with policy package. Red lines: Freeways (*Autobahn* and *Schnellstraße*), green lines: the Austrian railway network

## 5   Discussion

Our results highlight the need and potential for policy packages to increase public acceptance of restrictive policies, as found for example by [26]. While the introduction of a fuel price increase alone has disastrous effects on satisfaction and acceptance (as evidenced by actual fuel price increases worldwide in times of energy crisis), the addition of complementary measures can outweigh these effects. The most critical measure in this regard is the expansion of public transport, which can be understood as a prerequisite for all other measures (see also Fig. 4 in [23] for a graphical illustration). While there is still a very strong focus on urban areas, for example in [20], the challenges and relevance of transportation and the realities of people's lives outside of major cities should not be underestimated, as our findings make clear. This is also reflected in the comparatively higher acceptance of restrictive measures in urban areas where the infrastructure for public transport is already well-developed and there are more alternatives to the use of the car in contrast to rural areas.

While both satisfaction and acceptance would be highest if public transport was expanded exclusively, this does not make sense from a climate mitigation perspective. As [3] clearly state, restrictive measures are absolutely necessary to successfully reduce GHG emissions from transportation. Although the most supported packaging approach, namely inner-city driving bans and the expansion of public transport, already includes one type of restriction, it still does not take into account the need for

fundamental change not only in urban, but in all residential areas. In this respect, it is very interesting that the combination of a fuel price increase and the expansion of public transport leads to a higher acceptance than the combination of all measures. However, the combination of all measures is the only package with more than one push measure in place. Although not the most preferred policy package, including all four measures still results in acceptance levels over 60%, therefore indicating support.

The results of this study also have to be discussed critically in the context of policymaking. While, on the one hand, there is the greatest support for expanding public transportation, some regions and counties may find it difficult to create better connections and a good public transportation network due to time and financial constraints. In addition, solutions that work well in urban, high-density settings may not be translatable to the needs and realities of rural areas. Therefore, the implications of what a good public transportation system should provide for different residential areas need to be further elaborated. At the same time, we find that the often expressed fear of low acceptance of restrictive measures does not apply once such pull measures that increase the attractiveness of alternative modes of transportation are in place. This gives hope that a more diverse set of policy instruments finds use in passenger transport in the near future.

As with any study, this work comes with some limitations. As previously mentioned in Sect. 3, network effects are not rigorously tested but are used rather conservatively. Further elaborations on the role of social and personal needs (setting $\alpha$), and the richness of social interactions (clustering effects and assortativity regarding network topology as well as rural to urban differences) are necessary. However, a general problem that social simulations are currently facing is the lack of empirical data regarding opinion diffusion on networks, which will be hard to overcome for such a model as well. Another challenge is the quantification of indirect effects, e.g. effects on car travel satisfaction stemming from public transport policies, which is planned in a future version of the model.

We emphasize that the bottom-up approach to agent interactions allows the use of interaction topologies such as power-law and small-world networks. While these two topologies performed very well, random networks failed to reproduce the empirical data. Therefore, we argue that a mean-field approach or a top-down approach for this particular mathematical design would not match empirical results very well, as topological features were found to be essential.

## 6   Conclusion

The here presented ABM is used to evaluate policy package acceptance of push and pull measures. The model is informed and validated by empirical data on single policy measures, mode choice and environmental awareness of different residential areas. According to the results presented, we would like to emphasize that packages that include the expansion of public transport are highly appreciated in both rural

and urban areas. Combinations of push and pull price-related measures (fuel price increase and lower public transport costs) perform considerably worse. In general, the combination of push and pull measures can significantly increase acceptance for restrictive measures.

# References

1. Adnan, M., Outay, F., Ahmed, S., Brattich, E., Di Sabatino, S., Janssens, D.: Integrated agent-based microsimulation framework for examining impacts of mobility-oriented policies. Pers. Ubiquit. Comput. **25**(1), 205–217 (2021)
2. Ahanchian, M., Gregg, J.S., Tattini, J., Karlsson, K.B.: Analyzing effects of transport policies on travelers' rational behaviour for modal shift in Denmark. Case Stud. Transp. Policy **7**(4), 849–861 (2019)
3. Axsen, J., Plötz, P., Wolinetz, M.: Crafting strong, integrated policy mixes for deep CO2 mitigation in road transport. Nat. Clim. Change **10**(9), 809–818 (2020)
4. Banister, D., Hickman, R.: Transport futures: thinking the unthinkable. Transp. Policy **29**, 283–293 (2013). https://doi.org/10.1016/j.tranpol.2012.07.005
5. Drews, S., van den Bergh, J.C.: What explains public support for climate policies? A review of empirical and experimental studies. Clim. Policy **16**(7), 855–876 (2016)
6. Engler, D., Groh, E.D., Gutsche, G., Ziegler, A.: Acceptance of climate-oriented policy measures under the COVID-19 crisis: an empirical analysis for Germany. Clim. Policy **21**(10), 1281–1297 (2021)
7. Givoni, M.: Addressing transport policy challenges through policy-packaging. Transp. Res. Part A: Policy Prac. **60**, 1–8 (2014)
8. Hajinasab, B., Davidsson, P., Persson, J.A., Holmgren, J.: Towards an agent-based model of passenger transportation. In: International Workshop on Multi-Agent Systems and Agent-Based Simulation, pp. 132–145. Springer (2015)
9. Heinfellner, H., Ibesich, N., Lichtblau, G., Stranner, G., Svehla-Stix, S., Vogel, J., Wedler, M., Winter, R.: Sachstandsbericht mobilität und mögliche zielpfade zur erreichung der klimaziele 2050 mit dem zwischenziel 2030. UBA. Rep-0667. Im Auftrag des BMVIT (2018)
10. Helbing, D.: Agent-based modeling. In: Social Self-Organization, pp. 25–70. Springer (2012)
11. Huang, X., Lin, Y., Zhou, F., Lim, M.K., Chen, S.: Agent-based modelling for market acceptance of electric vehicles: evidence from China. Sustain. Prod. Consum. **28**, 206–217 (2021)
12. Huétink, F.J., van der Vooren, A., Alkemade, F.: Initial infrastructure development strategies for the transition to sustainable mobility. Technol. Forecast. Soc. Change **77**(8), 1270–1281 (2010)
13. Köhler, J., Whitmarsh, L., Nykvist, B., Schilperoord, M., Bergman, N., Haxeltine, A.: A transitions model for sustainable mobility. Ecol. Econ. **68**(12), 2985–2995 (2009)
14. Maggi, E., Vallino, E.: Understanding urban mobility and the impact of public policies: the role of the agent-based models. Res. Transp. Econ. **55**, 50–59 (2016)
15. Maggi, E., Vallino, E.: Price-based and motivation-based policies for sustainable urban commuting: an agent-based model. Research in Transportation Business & Management **39**, 100,588 (2021)
16. Mehdizadeh, M., Nordfjaern, T., Klöckner, C.: A systematic review of the agent-based modelling/simulation paradigm in mobility transition. Simul Paradigm Mobil Trans (2022)
17. Mueller, M.G., De Haan, P.: How much do incentives affect car purchase? Agent-based microsimulation of consumer choice of new cars-Part I: Model structure, simulation of bounded rationality, and model validation. Energy Policy **37**(3), 1072–1082 (2009)
18. Müller, M., Reutter, P.O.: Course change: navigating urban passenger transport toward sustainability through modal shift. Int. J. Sustain. Transp. 1–25 (2021)

19. Ormerod, P., Rosewell, B.: Validation and verification of agent-based models in the social sciences. In: International Workshop on Epistemological Aspects of Computer Simulation in the Social Sciences, pp. 130–140. Springer (2006)
20. Pörtner, H.O., Roberts, D.C., Adams, H., Adler, C., Aldunce, P., Ali, E., Begum, R.A., Betts, R., Kerr, R.B., Biesbroek, R., et al.: Climate change 2022: impacts, adaptation and vulnerability. In: IPCC Sixth Assessment Report (2022)
21. Silvia, C., Krause, R.M.: Assessing the impact of policy interventions on the adoption of plug-in electric vehicles: an agent-based model. Energy Policy **96**, 105–118 (2016)
22. Sopha, B.M., Klöckner, C.A., Febrianti, D.: Using agent-based modeling to explore policy options supporting adoption of natural gas vehicles in Indonesia. J. Environ. Psychol. **52**, 149–165 (2017)
23. Thaller, A., Posch, A., Dugan, A., Steininger, K.: How to design policy packages for sustainable transport: balancing disruptiveness and implement ability. Transp. Res. Part D: Transp. Environ. **91**, 102–714 (2021)
24. Tsoi, K.H., Loo, B.P., Banister, D.: "Mind the (policy-implementation) gap": transport decarbonisation policies and performances of leading global economies (1990–2018). Global Environ. Change **68**, 102–250 (2021)
25. Unruh, G.C.: Escaping carbon lock-in. Energy Policy 317–325 (2002)
26. Wicki, M., Huber, R.A., Bernauer, T.: Can policy-packaging increase public support for costly policies? Insights from a choice experiment on policies against vehicle emissions. J. Public Policy **40**(4), 599–625 (2020)
27. Wolf, I., Schröder, T., Neumann, J., de Haan, G.: Changing minds about electric cars: an empirically grounded agent-based modeling approach. Technol. Forecasting Soc. Change **94**, 269–285 (2015)

# Tools and Methods

# Formalising Agent Reasoning—The Paso Doble of Data and Theory

**Nanda Wijermans** and **Harko Verhagen**

**Abstract** One of the core assumptions made when building agent-based simulation models is how the agents decide or reason about the action to take next. The mode of reasoning is usually the same for all agents and over time within the simulation run. However, is this adequate? There exist several frameworks that describe multi-mode reasoning, however how do we know what we need? To engage with this core question, we reflect on this modelling process, by using CAFCA—one of these multi-mode frameworks—and reflect on the reasoning dimension in a social dilemma decision situation. More specifically, using existing qualitative inquiry on group dynamics in a common pool resource dilemma—not designed to elicit different types of reasoning—we introduce our hunt for reasoning hints and reflect on what insights/data we would need to make an informed decision about the reasoning(s) in our modelling and how to obtain this.

## 1 Introduction

One of the core assumptions made when building agent-based simulation models is how the agents decide or reason about the action to take next. The mode of reasoning is usually the same for all agents and over time within the simulation run. This however raises the question what type of reasoning is adequate and if it should or should not vary between agents and/or over time. Of course the answer to this question depends what the model is to represent and investigate and thus on the decision situations

N. Wijermans (✉)
Stockholm Resilience Centre, Stockholm University, Kräftriket 2b, 10691 Stockholm, Sweden
e-mail: nanda.wijermans@su.se

H. Verhagen
Department of Computer and Systems Sciences, Stockholm University, PO Box 7003, 16407 Kista, Sweden
e-mail: verhagen@dsv.su.se

that (may) occur within the simulation. For realistic modelling of human decision-making in more complex situations and over an extended period of time, we claim that different modes of reasoning should be taken into account and included in the model. There have been some suggestions for multi-reasoning mode frameworks such as Thinking fast and slow [4], the Consumat model [3], the Model Social Agent [1], and CAFCA [2]. We will discuss these in more detail in the next section, but for now, we detail the decision-making situation we will focus on.

A common pool resource dilemma, or more broadly social dilemma research, can reflect a wide range of situations, however, the take on reasoning is rather similar, i.e. only one type of reasoning: (bounded) rational actors [7], typically limited by being short-sighted and self-serving. Social dilemma or the study of the individual interest versus the collective benefits—is a vast research domain that involves theory, models and empirical data to understand, predict, and manage situations in which humans as a collective need to overcome their individualistic tendencies by working together and avoiding detrimental effect for everyone and—in the long run—benefit individually, such as situations in which over-fishing can occur, enjoying social health care, etc. Many problems humanity faces can be regarded as a social dilemma. Social dilemmas are typically framed and studied in one of the three typical stories/metaphors: the Prisoners Dilemma (2-player games), the problem of providing Public Goods (social fence—contribute first, benefit later) and The tragedy of the Commons (social traps—resist temptation, benefit for all) [5]. However, these three frames do not capture all social dilemmas that exist [5]. As much as this has structured the research, it also enabled the existence and persistence of blind spots [5]. We believe that the reasoning assumptions are one of them and we explore the way we can improve our self-reflection, to use and obtain empirical evidence for the (varying of) reasoning types.

In this paper we will investigate one of the multiple reasoning modes frameworks (CAFCA) and apply this to empirical data gathered (and lacking) in the development of an agent-based model of collective and sustainable common pool resource use (AgentEx). We will then introduce our self-reflection, which we hope to deepen and discuss during the conference, by reflecting on a decision situation that we have modelled as well as conducted empirical research on. Thus, the paper is about modelling practice and the interaction between modelling (rules) and gathering empirical (qualitative) data with theoretical models in mind to base the modelling on.

## 2   Multiple Reasoning Models Frameworks

Different areas of the social and behavioural sciences aim to understand behaviour at different levels, corresponding to different decision situations. The first two models are based on different strands of psychological theories and empirical research. In Kahneman's research (with various other researchers and co-authors, duly presented in [4]) into various types of biases and other deviations from rational decision making resulted in two general decision making modes—fast (system 1) and slow

(system 2) thinking. System 1 is suggested to cover 98% of all thinking, lacking any self-awareness or control and covers automatic or unconscious "thinking" including heuristics. System 2 supposedly covers the remaining 2% where deliberation and rational thinking and information gathering are taking place. The named percentages do, however, not reflect the real world as such, but the data found in the experiments devised by Kahneman c.s. The Consumat model [3], which also has its base in psychological theories and models of decision making, is more varied as it consists of four different modes. One of these consists of automatic decision-making, in which one could see the other 3 modes as varieties within system 2. It differs from Kahneman's model as it has clear mechanisms on why and how the other modes come into play. In short, the uncertainty regarding the outcome of a decision is based on previous experiences. Including observations of results of decisions other agents have made (adding a social perspective largely lacking in Kahneman's model) as well as what level constitutes satisficing in the tradition of [7]. The Consumat authors later added social networks to their model to extend the social part of the decisions under consideration—consumer decision situations regarding the purchase of goods. The original Consumat model distinguishes four different reasoning modes while the extended versions consists of six such modes.

The remaining two frameworks are more extensive and build upon analytical categories with a range of social science based aspects. The Model Social Agent, as developed by Carley and Newell [1] is an ambitious suggestion for such a framework. They distinguish two dimensions along which agents can be categorised: the information processing capabilities of the agent sand the knowledge the agent possesses. These dimensions are also used to categorise theories on the behaviour and decision making of artificial and human agents. In this framework, social science concepts characterise different levels of sociality while the agent-internal frameworks regarding information processing come mainly from the cognitive sciences. The simplest and least human-like agent one can imagine using these dimensions is an agent who is omnipotent (i.e., without any limitations at all on its information processing capabilities) possessing only knowledge about the task it is working on. At the other extreme, we find an agent who is both emotional and cognitive, with knowledge of the task, as well as other agents, interacting in real-time within a social structure with social goals and cultural and historical knowledge of that social structure. In total, the framework forms a six (knowledge involved) by five (information processing) matrix, resulting in thirty reasoning modes. The Model Social Agent framework is particularly strong in its encompassing of subtly differing social theories. The Consumat and the Model Social Agent frameworks inspired the development of the Contextual Action Framework for Computational Agents (CAFCA) [2] which we will use in this paper. CAFCA distinguishes two dimensions of decision-making context that together frame what models of reasoning can be applied in the resulting context as it is seen by the agent, see Fig. 1. One dimension describes the type of reasoning: habitual, strategic (goal-driven), or normative. The other dimension pinpoints the level of sociality: individual, social, or collective. In the individual mode the agent interprets the decision as independent of others. In the social mode agents recognise other agents in the situation but sees oneself as distinct from or in competition with

**Fig. 1** Contextual action framework for computational agents (CAFCA) applied to common theories. Adapted version of [2]

| | | SOCIALITY DIMENSION | | |
|---|---|---|---|---|
| | | INDIVIDUAL | SOCIAL | COLLECTIVE |
| REASONING DIMENSION | HABITUAL | Repetition | Imitation | Joining-in |
| | STRATEGIC | Rational choice | Game theory | Team reasoning |
| | NORMATIVE | (institutional) rules | (social) norms | (moral) values |

them. In the collective mode the agent does not only recognise others but perceives itself as belonging with the others, as a member of a collective or team. Together, these form a matrix of nine reasoning modes.

## 3 Learning to Paso Doble: Theory and Data Back and Forth, Step 1

The first author engages within a long-term collaboration -AgentEx -in which agent-based modelling and controlled behavioural experiments are combined to contribute to the understanding of the sustainable and collective use of natural resources [9]. The research consists of developing and testing explanations (experimental outcome → ABM as explanation able to reproduce patterns) and exploration, designing data collection around an behavioural experiment to be able to generate explanations and hypotheses on for instance the role of perception on behaviour [8]. The decision context of both the behavioural and simulation experiments concern a social/common pool resource dilemma. While in the past we have reflected on the sociality dimension of AgentEx [10], here, with use of the CAFCA, we will reflect on the reasoning dimension. Given the data collection around an experiment, we were curious what we would learn from empirical insights on group dynamics found in the empirical data of a study on the decision context of AgentEx, which was not designed to extract reasoning information. This analysis has two purposes: (1) to investigate what can we learn about reasoning, and (2) to formulate/design the empirical inquiry needed to obtain relevant reasoning information.

## 3.1 Description of the Empirics—Data Collection

To obtain more information about the group dynamics for an ABM of sustainable and collective common-pool resource use, the first author conducted interviews with the experiment team members that performed controlled behavioural experiments. It concerns a framed field experiment in the form of a dynamic common-pool resource (CPR) game designed to capture behavioural responses of resource-dependent small-scale fishers to potential resource scarcity [6]. As described in [9] this behavioural experiment is a so-called'pen-and-paper experiment' in which 4 participants (small-scale fishers in our case) sit at a table, get information on paper, and are accompanied by an experiment team who guide them through the game rounds. The group plays for the duration of 14 rounds, a duration that is not known to the participants. During each round 4 participants could: (i) communicate face-to-face, (ii) individually and anonymously harvest resources by writing down how much resource units they want, and then (iii) were informed by the experiment leader about the resulting fish-stock (after harvesting and renewal of the resource). Based on how they played they received payment for each unit of resource after the game. This set-up is common for behavioural experiments on CPR use.

To explore the reasoning that occurs in these types of decision context empirically, we qualitatively examined the group dynamics processes in this behavioural experiment. More specifically, we used the 6 debriefing interviews that were conducted with the experiment teams after the behavioural experiment with Thai fishers in different coastal communities. The purpose of these interviews were to get a feeling for the group dynamics to support the formalisation process. In total there were 42 behavioural experiments, the choice for which group would be interviewed depended on the availability of the experiment team (afternoon session only) and the presence of the first author (first half of the field experiments). The aim of the project for this data collection is to formalise the influence of perception of change in the resource on their actions via their internal characteristics and processes. We asked questions about the group dynamics, their perceptions and attributions, but also on whether they felt like a group, whether this changed throughout the experiment etc.

## 3.2 Looking for Clues

To investigate what can we learn about reasoning in this study, the first author re-analysed the six interviews. Through listening to the recordings and reading the written summaries hints of reasoning were identified with the CAFCA reasoning levels in mind: habitual, strategic and normative reasoning. We reflected on the reasoning for each choice in the decision situation of the experimental design concerning: (1) whether to communicate or not, and (2) how much to take out. The first decision occurs during the communication stage: while communicating each participant decides whether to communicate and if then about what to communicate: (a) the

amount the group should take with reasons, or (b) questions to help you understand the dynamics better and/or (c) relational information to assess whether you trust the others. The second decision relates to the amount one will harvest being based on a potential agreement made OR by following one's own individual reasoning. While acknowledging that it is hard to derive the reasoning from each individual participant, the communication stage allowed for detecting hints about potential group reasoning and the actual behaviour in combination with the overall group dynamics sometimes provide hints about the reasoning behind what to communicate and deciding on the harvest level. In total five out of the six debrief interviews signals/hints about the reasoning in the group and/or the individual participants could be identified.

**Group 1: making it happen together | strategic, habitual and normative reasoning.** Core characteristic of group 1 is that they manage the resource well. All participants engage in communication, and while one participant does not really "get it", together they make it work. Their reasoning seems focused on taking the optimal amount, which can be interpreted as the outcome of strategic reasoning. However, over time, more and more resembles habitual decision making. It seems that the participants find their strategy and stick with it, while the state of the resource reinforces that everyone is on board.

On the individual agent level, participant #1.3 seems to set the amount for the agreement that is followed, making us assume this person reasons strategically and follows that themselves when extracting (decision 1 and 2). For participant #1.4, who does not seem to understand the required reasoning, the reasoning seems to be more in following the agreement when performing the extraction (decision 2). One could see this as normative level reasoning ("following the rule"). For the remaining two participants, who both actively contributed by sharing their understanding of the resource, the behaviour could indicate they think strategically in contributing and reaching an agreement (decision 1). However, when actually performing the extraction, it could well be that they follow the agreement rather than their own reasoning (decision 2). While we can impossibly know for sure, with person #1.3 and #1.4 we feel there are some reasoning hints, while for #1.1 and #1.2 we cannot determine their reasoning for either decision.

**Group 2: the one that ruins it | strategic reasoning.** Core characteristic in group 2 is that while three persons communicate, one participant (#2.3) remains mostly silent, but seems to pretend to not understand how the system works. The other three communicating participants align, understand the resource, what is optimal, and agree on what to do to reach the optimal stock size and regeneration rate. This subgroup seems to reason strategically. Throughout the game too many resources are taken, the three try to convince #2.3 to take less, but #2.3 does not listen. When the three notice that something was off with #2.3, one participant increased the extraction after while and one decreased the extraction to compensate. In the end, one of the participants extracts the total remaining stock.

On the individual level, each of the three participants shows an understanding of the system and a persistent focus on the optimal points to a dominant strategic interaction. However, over time, due to the bad performance, more participants change

their behaviour (take less and compensate, take more as things are lost anyway). The changes are however along CAFCAs sociality dimension, while the reasoning remains strategic. Similarly, person #2.3 seems to be strategic in not communicating (decision 1) and harvesting following its own reasoning (decision 2) that could be strategic every round or habitual, there is no way to know.

**Group 3: silence is not golden | normative and strategic reasoning.** Core characteristic in group 3 was that 3 participants do not really communicate, only participant #3.4 tries to communicate and convince everyone what to do, as he understands very well how the system works. At the same time #3.4 is the youngest of the group and in the Thai cultural context, the younger ones listen to the older ones. However the oldest ones do not understand the resource and remain silent. The group depleted the resource prematurely.

On the individual level, the only cues we explicitly get are from participant #3.4, who seems to follow strategic reasoning when deciding how much to harvest, while at the same time bending the normative 'rules of engagement' regarding to be heard or not. The participant tries to convince the others of the best course of action that allows them all to have income. Participants #3.2 and #3.3 seem more concerned with what is appropriate behaviour in this context given the age differences, which would point to normative reasoning, and busy with to whom they should be listening to, i.e., they silently consider whether they need to stick to the same cultural rules or renegotiate the cultural agreement (decision 2).

**Group 4: conflict and reconciliation, damaging trust | strategic and normative reasoning.** Core characteristic: a turbulent group due to changes in who is influencing the group agreement, a temporary stop of communication occurs after massive dip in the resource and a conflict. In the end the group turns things around into cooperation, however it too late to maintain the resource. In the beginning, everyone followed participant #4.4, who was always taking more than agreed. After the conflict everyone followed the suggestions of participant #4.2. There were more tensions, as participant #4.4 wanted to end the game to take time to pray (religion), whereas participant #4.2 wanted to play as long as possible. Also, participant #4.2 told participant #4.4 to "Don't be greedy" because participant #4.4 always said he wanted to catch a lot.

In terms of reasoning the communication in the group was about getting to influence the others to follow them. It seems as if participant #4.4 was strategical changing the norm, encouraging other into normative reasoning, to then take advantage of that. Overall the group seemed to predominantly engage in normative reasoning and following the strategic suggestions of participant #4.4 and later participant #4.2.

**Group 5: one that spoils it all | strategic and ..?** Core characteristic of this group is that they are a blend of people who partially do and partially do not understand the resource dynamics as they deplete the resources rather quickly. Some participants focus on talking and influencing others while other participants just keep to themselves. Being young also affects the game here, as the younger participants have no influence and seem to lose confidence after a big loss in biomass. There is no explicit agreement on amount of harvest, just that they want to be at the optimal rate (according to participant #5.1).

On the individual level, based on the individual interviews, participant #5.4 seems to really understand the system and wanted to be at the optimal (strategic, decision 2), shared this but due to their age had no influence and further refrained from communicating (normative, decision 1). Participant #5.1 understands the system as well and also tries to convince the others, which everyone seems to accept and agree with. When looking at the extraction levels, it shows participant #5.2 is always taking more, the experimental team indicates the person does not understand the system, however this participant takes even more in phase 2 which makes one wonder how non-strategic this is. And if this participant was not reasoning at the strategic level nor the normative level, what reasoning was there? Were the choices just random? It was definitely not habitual either...

## 4   Discussion

Our analysis shows that reasoning modes differ between agents and over time for the same agent, despite the relatively short duration of the experiment and the mildly complex system dynamics (akin to the occupational system dynamics of the participants) and the small group size. The analysis also points out how hard it is to characterise reasoning modes from the (sparse) interview data we have. Mapping the data to the CAFCA reasoning modes or even levels is not straightforward, even if we see different modes appear. The reason for these differences are opaque and even more so for the mode change for one agent where the how, when, and why of mode switching is an even harder challenge. For example, is it only a positive reinforcing situation that leads to habitual behaviour, like we saw in group 1? And that a negative reinforcement leads to change, be it in reasoning or on the sociality dimension? Finally, interpreting the communication during the experiment is, again, not easy. For instance, it is not always clear if the reported reasoning during the experiment is actually how the participant reasons or if this is meant to 'manipulate' others to behave in a certain way.

This then leads to the holy grail of data gathering, what kind of questions and observations would allow us to infer the deeper meaning of what is going on. Other research that elicits reasoning, such as expert system development, often focuses on deliberate reasoning and uses the 'talk-aloud', 'walk-along' methodologies of knowledge elicitation. This leaves open the issue of how to track subconscious processes leading to decisions. Perhaps a measure for this could be the effort or amount of time spend (or available and spend) which indicates if reasoning is deliberative or habitual, as both the Consumat model and Kahneman's two systems model point to this. We could also turn to more anthropological ways of data gathering to aim for "thick descriptions" that help us explain in a deeper way what is going on. This would put higher demands on the design and execution of the data collection.

Thus, we have a list of open issues that we are looking forward to discuss.

- What are potential reasons for mode switching,
- Starting from our data, how would other researchers map these on CAFCA matrix,
- What data would be relevant to collect (and how) to get to reasoning?, and
- What theories are relevant yet missing from CAFCA or other frameworks?

## References

1. Carley, K., Newell, A.: The nature of the social agent. J. Math. Soc. **19**(4), 221–262 (1994)
2. Elsenbroich, C., Verhagen, H.: The simplicity of complex agents: a contextual action framework for computational agents. Mind Soc. **15**(1), 131–143 (2016)
3. Jansen, M., Jager, W.: An integrated approach to simulating behavioural processes: a case study of the lock-in of consumption patterns. J. Artif. Soc. Soc. Simul. **2**(2) (2000)
4. Kahneman, D.: Thinking, Fast and Slow. Farrar, Straus & Giroux (2011)
5. Kollock, P.: Social dilemmas: the anatomy of cooperation. Ann. Rev. Soc. **24**(1), 183–214 (1998)
6. Lindahl, T., Schill, C., Jarungrattanapong, R.: Beijer discussion paper 276: the role of resource dependency for sharing increasingly scarce resources: evidence from a behavioural experiment with small-scale fishers. In: Beijer Discussion Paper Series (2021)
7. Simon, H.A.: A behavioural model of rational choice. Q. J. Econ. **69**(1), 99–118 (1955)
8. Wijermans, N., Schill, C., Lindahl, T., Schlüter, M.: AgentEx (2016). https://doi.org/10.25937/js95-6d78
9. Wijermans, N., Schill, C., Lindahl, T., Schlüter, M.: Combining approaches: looking behind the scenes of integrating multiple types of evidence from controlled behavioural experiments through agent-based modelling. Int. J. Soc. Res. Methodol. **4**, 1–13 (2022). https://doi.org/10.1080/13645579.2022.2050120
10. Wijermans, N., Verhagen, H.: Fishing together? - exploring the murky waters of sociality. In: Dam, K.H.V., Verstaevel, N. (Eds.). Multi-Agent-Based Simulation XXII - 22nd International Workshop, MABS 2021, Virtual Event, May 3-7, 2021, Revised Selected Papers. Lecture Notes in Computer Science, vol. 13128, pp. 180–193. Springer (2021). https://doi.org/10.1007/978-3-030-94548-0_14, https://doi.org/10.1007/978-3-030-94548-0_14

# Model Mechanisms and Behavioral Attractors

**H. Van Dyke Parunak** ⓘ

**Abstract**  In social modeling, a *computational environment* runs a *model* that represents the *world*. The states the model explores (its *behavioral attractor*) are typically fewer than its description suggests. The mapping between model and attractor depends not only on its *parameters* (exploring variants of the world) and its *conventions* (imposed by the computing environment), but also its *mechanisms* (components of the model representing selected dimensions of the world). We illustrate the impact of different mechanisms on the attractor. In our case, in general, the more mechanisms one implements, the smaller the attractor ("the more you model, the less you see"), but with unexpected twists.

**Keywords**  Agent based modeling · Social simulation · Complex dynamics · Model parameters · Model mechanisms · Behavioral attractor

## 1  Introduction

The user of a social model expects the model to generate a range of behaviors. For example, how many distinct behaviors can the actors manifest? How does their spatial distribution vary over time? The range of behaviors generated by a running model (the system's *behavioral attractor*) is usually smaller than the static model suggests.

The mapping between a model and its attractor can depend on three different sets of variables: *parameters*, *conventions*, and *mechanisms*. Each of these describes a different component of the modeling enterprise, in which a *computational environment* runs a *model* that represents the *world* (Fig. 1).

- *Parameters* describe the world that the model represents. Varying them explores how the world might behave if its characteristics (e.g., relative group sizes) change.
- *Conventions* are unrelated to the real world but imposed by the computational environment (such as agent execution order on a von Neumann machine or agent

H. Van Dyke Parunak (✉)
Parallax Advanced Research, Beavercreek, OH 45431, USA
e-mail: van.parunak@gmail.com

**Fig. 1** Central relations



behavior at arena boundaries), and varying them explores the degree to which the behavioral attractor is an artifact of that computational environment.

- *Mechanisms* are model components that reflect facets of the world. For example, real social actors have short term preferences and strategic goals that guide their choices, subject to constraints among available options and the actions of other actors. Not every social model has a mechanism for each of these (preferences, goals, option constraints, interactions), and no social model has a mechanism for every possible dimension.

Modelers assume that a model with fewer mechanisms than the world's facets can still give useful information. Most modeling frameworks offer few alternative mechanisms, seducing modelers to ignore the impact of mechanism choice. SCAMP (Social Causality using Agents with Multiple Perspectives) [13], a causal language and simulator for social scenarios, has a rich array of mechanisms that can be activated independently of one another. In general, the more mechanisms we activate, the smaller the attractor ("the more you model, the less you see"), but interactions among mechanisms lead to anomalies. For instance, a more constrained attractor may lie partly outside less constrained ones with the same conventions and parameters. Adding mechanisms can not only sharpen the model's focus, but also shift its location.

These results are of immediate interest to teams who are using SCAMP. In addition, our methods should be helpful to other modelers in understanding the implications of their choice of mechanisms. Our exploration of SCAMP's dynamics is a concrete example of what might be done in other frameworks.

Section 2 summarizes related work. Section 3 describes the mechanisms of SCAMP that these experiments vary. Section 4 describes our methodology. Section 5 presents the experimental results. Section 6 discusses their implication for interpreting the results of a SCAMP run, highlights implications of this experiment for other social modeling systems, and outlines future work.

## 2 Related Work

We expect behavior to vary with model *parameters*, which are the focus of most studies of the dynamics of agent-based systems (e.g., [2–4, 20]), including studies of tipping points (parameter values where behavior changes discontinuously, leading to a phase shift) and lever points (parameters whose change has a lasting, directed effect) [1, 15]. Wolfram [19] identifies four distinct classes of one-dimensional 0–1 nearest-neighbor cellular automata, varying only the update rule, the key model parameter. Verification methods such as sensitivity analysis [5] (p. 24) or comparison of agent trajectories with observed data also explore behavioral changes when parameters change, but not the impact of changing conventions or mechanisms.

Studies of the impact of *conventions* imposed by the computing environment are less common, but revealing. For example, a differential equation model and an agent-based model can yield qualitatively different results for the same parameters [16, 18]. Restricting ourselves to agent-based modeling, on a von Neumann machine, agents can run only one at a time, and different scheduling disciplines for entities that in reality execute concurrently have repeatedly been shown to lead to different results [6, 8]. Núñez-Corrales [9] reviews the extensive literature on the impact of scheduler synchrony.

This study focuses neither on the parameters that vary the world explored by a model nor on the conventions imposed by computation, but on differing sets of mechanisms that the model uses to represent facets of the world. Naively, one hopes that even a primitive model will be useful, and that adding more mechanisms will add more detail to the results of the initial model. Unexpectedly, such refinements can also move the focus, and cause other anomalies. We know of no other ABM work that demonstrates this effect, because most modeling frameworks do not offer multiple mechanisms that can be activated independently of one another.

## 3 The SCAMP Causal Modeling System

This section explains enough of SCAMP's structure to motivate our experiments. For further details, see [11, 13]. Our experiments use two of SCAMP's four perspectives.

1. A *causal event graph* or CEG is a directed graph whose nodes represent types of events in which agents can participate.
2. A *hierarchical goal network* or HGN is a directed acyclic graph that models the goals of a group of agents and how those goals are related to the levels of participation on events in the CEG. Leaf nodes in the HGN are linked, or *zipped,* to event nodes that either support or block them.

The CEG has two sorts of edges:

1. An *agency edge* from node A to node B means that an agent currently participating in an event of type A may consider an event of type B as its next activity. A chain of

agency edges defines a plausible *narrative* of the agent's experience. Depending on its group membership, an agent has *agency* for a subset of the nodes in the CEG, and can move between two nodes only if it has agency in both of them. Most nodes have multiple successors, making the CEG a *narrative space* [14] that captures many possible narratives. The main output from a SCAMP model is the history experienced by each agent. Agency edges are obligatory.

2. Sometimes one event causally constrains another even though no agent has agency for both events. For example, an act of God such as a pandemic may hinder events in which people gather together, or enhance hospitalization events. SCAMP captures these relationships with *influence edges*. Influence edges are optional.

When an agent completes one event in the CEG, it selects the next based on two vectors. Each event has a *feature vector* that describes the event's effect on agent wellbeing, how urgent the event is to satisfying the HGNs to whose leaf goals it is zipped, and how extensively agents of each group have participated in it recently. Each agent carries a *preference vector* over the same space. To choose its next event, the agent

1. computes the dot product of its preference vector and the feature vector of each accessible event type in the CEG,
2. exponentiates each dot product so that it is non-negative, defining a roulette,
3. adjusts the presence and size of segments with incoming influence edges based on the participation levels on events at the origins of those edges,
4. and spins the roulette.

In step 3, a *prevent* or *enable* influence edge can remove or add an event to the roulette that guides agent choice, changing the structure of the CEG dynamically as participation levels on influencing events vary.

Each agent adjusts its roulette before spinning by raising the size of each sector to its personal *determinism* level, modeling human deviation from pure rationality. An agent with determinism 0 makes completely random choices, while determinism 100 models a utility optimizer. Our experiments set agent determinism to 100, while our baselines set it to 0 to generate a random walk.

SCAMP uses polyagents [10], which represent each domain entity by a single *avatar* that can deploy a swarm of *ghosts*. The ghosts explore their avatar's possible next choices by looking ahead a fixed distance (here, two events). At each step, they form a roulette over all nodes in the CEG that are immediate successors to their current node, choose one node, and increment the node's presence feature for their group proportional to the value of the position in which they find themselves. The avatar chooses its next step by choosing probabilistically based on the presence features deposited by its ghosts. This mechanism simulates the well-documented psychological process of evaluating actions by mental simulation of possible outcomes [7].

We base our experiments on a model of civil strife inspired by recent history in Syria. The CEG in this model includes 460 event nodes with 1106 agency edges and 400 influence edges. The six HGNs, one for each group, include 122 goals or subgoals. 77 leaf goals are zipped to 177 event nodes.

## 4 Experimental Methodology

Our methodology has three parts.

1. Define how to measure *the space of behaviors.*
2. Identify the *mechanisms* that an instance of the model supports. A given set of mechanisms defines a *configuration.* We are interested in how the size of behavior space varies with the configuration.
3. Identify a configuration to represent an unconstrained *baseline.*

### 4.1 Defining Behavior Space

An analyst constructing a SCAMP model starts with the CEG, defining types of events that might occur in the domain and linking them into reasonable narratives for agents belonging to different groups. One useful measure of behavior space is how many of these event types the system actually explores. Two levels of exploration are meaningful. The first counts node *coverage*, in several ways:

1. How many nodes do *ghosts* visit in evaluating possible futures for their avatars?
2. How many nodes do ghosts *consider* in evaluating possible futures for their avatars? These are successor nodes to those nodes that the ghosts actually visit.
3. How many nodes do the *avatars* visit in carrying out their decisions?

We measure these values for multiple runs of each configuration, with different random seeds. In this paper, we run at least six runs per configuration.

We also look at how similar the sets of nodes under each measure are for repeated runs with different random seeds. Let $Q$ and $R$ be the sets of nodes explored (under one of the options above) for two runs of the same configuration, and let $S$ be the union of the sets explored by both runs. Then the *overlap* between $Q$ and $R$ is $|Q \cap R|/|S|$.

We hypothesize that as we add mechanisms, the numbers of distinct nodes in each category will drop (the attractors will shrink) while the overlaps will increase, because the system will be attracted into the same region of state space. As we will see, the data hold some surprises that yield important insight into the system's behavior.

These measures do not by any means exhaust those that could be considered. Since a commonly used output of SCAMP is the set of behavioral trajectories followed by the agents, one very relevant measure is the number of distinct trajectories that avatars execute. We leave that analysis for future work.

## 4.2   SCAMP's Mechanisms

SCAMP offers several mechanisms to capture different dimensions of the world.

The most basic is the *structure of the agency edges* in the CEG, which record the meaningful behavioral trajectories available to agents. Even if agents execute random walks, the branching factors differ along different paths, so that nodes only accessible along highly branched paths will have a lower probability of being sampled in a run of a given length than those with less ramified approaches.

The mean node degree restricted to agency edges in our example CEG is 4.74, not much more than the limit of 4 for an infinite square lattice, but degree in the CEG is highly variable. Consider a synthetic baseline of 460 integers randomly selected from [3, 6]. The mean is 4.5, comparable to our data, but Pearson's kurtosis for this synthetic baseline is 1.64, well below the threshold of 3 associated with normally distributed data. For our CEG, the kurtosis of node degree is 8.7, reflecting the heavy tail of nodes with high degree (up to a maximum degree of 21).

For comparison, we construct a rectangular directed lattice of 21 * 22 = 462 nodes, over which we do a random walk (with both ghost and avatar determinism set to 0). A random walk on a regular lattice with restart will visit every node, if it runs long enough. We expect the CEG to perform similarly. We also do a random walk over the CEG model itself, augmented with a single START and a single STOP node.

Psychological preference is modeled by the *feature space* that defines agent preferences and event features. Without preferences, ghosts perform a random walk in laying down the presence features that guide avatars. With preferences, ghosts will favor some nodes over others, based on the features that the model builder has defined for those nodes. We expect (a) agents using preferences will explore fewer nodes than those walking randomly, (b) overlap across runs will be greater with preferences than without, and (c) the longer the model runs, the more nodes will be visited.

SCAMP's HGNs model strategic reasoning. Each HGN monitors the recent participation level on event types to which it is zipped to assess its current *satisfaction*, then computes the *urgency* feature of each of these events. Agents respond to urgency according to their preferences. If an agent is running without preferences, the HGN is irrelevant. But if preferences are active, we expect HGNs to focus the agents' attention, reducing the number of nodes explored and increasing their overlap.

*Influence edges* model causal influences among event types between which agents do not move directly, modulating the probability of destination nodes dynamically based on participation levels on source nodes. Again, including this mechanism should reduce the number of nodes visited and increase their overlap.

**Fig. 2** Configuration lattice:
001 = influences, 010 =
HGN, 100 = preferences



A *configuration* is a binary string indicating active mechanisms. The first position shows whether (1) or not (0) preferences are active. The second position shows HGNs, and the third, the use of influence edges. Thus in 000, the only mechanism is the structure of agency edges, 100 indicates the use of preferences alone, 110 adds HGNs, and 001 is the use of influence edges alone. The decimal values of these strings identify configurations 0 (no mechanisms active) to 7 (all mechanisms active). Configurations 2 and 3 (HGNs without preferences) violate the assumptions of the model and are not included. Configurations 4–7 include preferences, configurations 6 and 7 include HGNs, and odd configurations include influence edges. Our configurations thus form a partial lattice (Fig. 2). All configurations use the same parameters to describe the world and run with the same conventions.

## 4.3  Random Baseline

In addition to a space within which the attractor is defined and mechanisms that might impact the attractor, we need a baseline against which to compare their impacts. We provide two baselines, L (the 21 * 22 lattice) and R (the CEG), with both ghost and avatar determinism set to 0 so that they ignore the roulette entirely. In configuration 0, unlike R, avatars have determinism 100, and follow their (randomly moving) ghosts.

## 5  Results

Our experiments illustrate how studying the behavioral attractor as a function of model mechanisms can yield valuable insights that confirm or correct our intuitions and call attention to behaviors that invite further study. Our study is exploratory, and we present most results as boxplots [17].[1] In some cases, we compute the significance of pairs of results using the one-sided Mann–Whitney U test. p-values greater than 0.05 are reported as not significant.

---

[1] The box extends from the upper to the lower quartile of a data series. The bold line marks the median. The whiskers extend to the most distant data points within 1.5 times the inter-quartile range of the quartile limits, and circles mark outliers. Comparing the inter-quartile boxes for two series is a good heuristic for whether they are the same or different.

We begin with summary plots that characterize the data and show the impact of run length on our measures. Then we examine how visits and overlaps vary with configuration, to see how adding mechanisms affects agent activity. Finally, we offer some summary statistics on the impact on our metrics of the three mechanisms we are studying: preferences, HGNs, and influence edges.
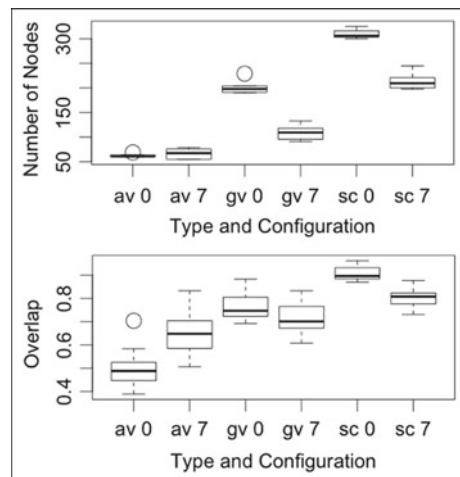
## 5.1   Making Friends with the Data

First, compare the coverage and overlaps (Fig. 3) of each measure [avatar visits (av), ghost visits (gv), and successors considered by ghosts (sc)] for the baseline configuration (000 ~ 0) and the most constrained (111 ~ 7). Ghosts visit fewer nodes than they consider, and avatars visit only a small fraction of those explored by ghosts. Added mechanisms reduce the number of nodes that the ghosts consider and visit, as expected, but the number of nodes visited by avatars is unchanged. Additional mechanisms focus the ghosts' attention more closely, but however broadly or narrowly the ghosts explore, an avatar chooses one path from those explored by its ghosts, and in a run of fixed length visits only a limited number of nodes. The avatar nodes are not the same in the two configurations, but by the structure of the program the coverage is the same size.

We expect overlap to increase with mechanisms, as agents focus their attention on fewer nodes. Figure 3 confirms this intuition for avatar visits, but overlaps for ghost visits and successors actually *decrease*, a phenomenon we discuss in Sect. 5.3.

In a regular directed lattice, coverage increases with run length. Most of our results are runs of 1000 Repast ticks. Figure 4 shows the effect of increasing run length to

**Fig. 3** Nodes visited (left) and overlaps (right) by types and configurations

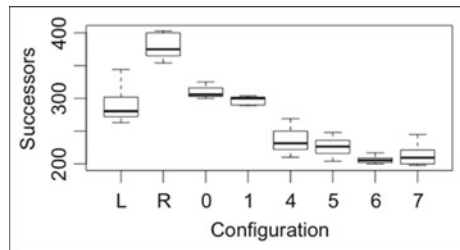**Fig. 4** Effect of run length on coverage



2000. We compare configuration 0 with 4, which (we will see) is particularly influential. In x-axis labels, the first digit (0, 4) is configuration, and the second (1, 2) is run length in k-ticks. The intuition is correct for avatar visits, and for ghost visits and successors in configuration 0. But for configuration 4, the preference mechanism leads the system to converge, and longer runs do not increase ghost visits or successors.

The median value of sc01, 306, leaves 162 event types in a typical run that the ghosts never consider. However, these 162 CEG nodes are not the same in each run. The median overlap is about 90%, and many runs show lower overlaps between pairs of runs for each configuration. For example, while the maximum successor nodes in any single run of configuration 0 is 325, all of the runs together explore 343 nodes. This still misses 117 nodes of the complete CEG, but suggests that multiple runs are at least as important as run length in sampling the causal graph adequately.

## 5.2 Impact of Adding Mechanisms

Figure 3 shows a clear reduction in coverage for successors and ghost visits between runs with no mechanism except the CEG, and all mechanisms. Figure 5 shows intermediate configurations. Ghost visits show the same pattern.

**Fig. 5** Successors by configuration

In the baselines, random walk on a lattice (configuration L) offers fewer successors to consider (and thus fewer ghost visits) than on the CEG (configuration R), reflecting the long tail in our model's degree distribution. Successors and ghost visits on baseline R are higher than with any mechanisms, which is not surprising.

As expected, both measures tend to decrease as we add mechanisms. Two details are particularly interesting.

1. The drop from configurations 0 and 1 (no preferences) to configurations 4–7 (with preferences active) is particularly large, suggesting that preferences have more influence on the system than do HGNs or influence edges.
2. Configuration 6 (preferences and HGNs without influence edges) appears to be *lower* than the more highly constrained configuration 7 (which adds the influence edges). This unexpected result shows an unanticipated but realistic interaction between the two mechanisms. An agent's goals (in life and in SCAMP) guide its actions by identifying high-priority events in which the agent should participate, and the usefulness of goals will decrease if other events block access to those urgent events through influence edges.

In contrast to successors and ghost visits, avatar visits do *not* change with more mechanisms. This observation is consistent with Fig. 3: while ghosts can explore more or less narrowly, each avatar follows only one path, and thus visits only a relatively constant number of nodes for runs of a given length.
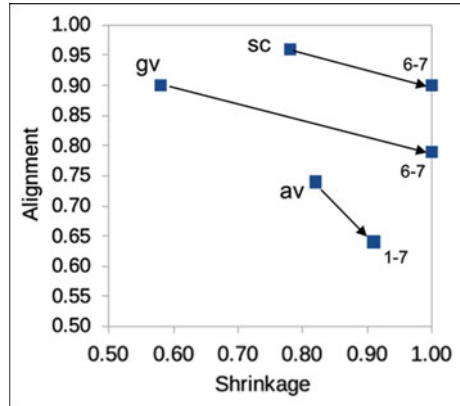
Adding constraints not only decreases attractor size (for gv and sc), but also shifts its location. Define the *alignment* of one configuration with another with a subset of its mechanisms as the percent of events in its attractor with that of the less constrained configuration, and *shrinkage* as the ratio of the size of the more constrained attractor to the less constrained. Figure 6 shows the alignment and shrinkage for six pairs of configurations, two for each metric. In each case, one pair compares the attractor for configuration 0 with that for configuration 7 ($0 \rightarrow 7$), while the other compares configuration 7 with a less constrained configuration ($6 \rightarrow 7$ or $1 \rightarrow 7$). As expected, shrinkage is greatest for the greatest increase in constraints ($0 \rightarrow 7$). But contrary to expectation, increasing constraints *reduces* the alignment between attractors. This result challenges the common assumption that ignoring facets of the real world gives a fuzzier but still essentially correct outcome. In fact, adding these facets can shift the model's output.

## 5.3 A Closer Look at Overlap

In addition to monitoring the coverage of nodes considered or visited (our approximation of a model's attractor), it is also useful to study the variation among the sets of nodes visited in different runs of the same configuration. Intuitively, we expect overlap to increase with number of mechanisms. This intuition must be qualified.

With significance $p = 2E-16$, avatar visits have lower overlap than ghost visits, and ghost visits have lower overlap than successor coverage. We hypothesize that

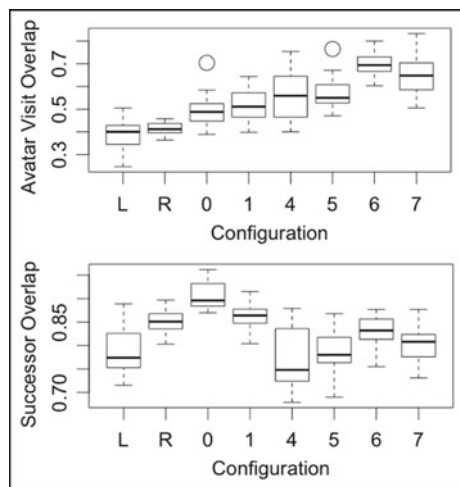**Fig. 6** Alignment of Successive Constraints: square = av, dot = gv, triangle = sc

this difference reflects the fact (Fig. 3) that there are far more successors than ghost visits, and far more ghost visits than avatar visits, out of a fixed number of nodes. Higher coverage of the CEG leaves fewer nodes on which runs can differ with each other.

Figure 7 shows how avatar and successor overlaps vary with configuration. Avatar overlap satisfies our intuition that with more mechanisms guiding agents into similar regions of the CEG, overlap should increase. Consistent with this dynamic, the baseline configurations L and R, with both ghosts and avatars executing random walks, have the lowest overlaps. Configuration 6 yields the highest overlap. Adding influence edges in configuration 7 reduces overlap, reflecting their interaction with HGNs.

Ghost overlap (not shown) is less intuitive. Overlaps between the sets of nodes visited by ghosts in different runs of the same configuration are invariant with configuration, and do not significantly differ from the baselines.



**Fig. 7** Overlaps by configuration

Overlaps in the successor metric are even more complex. Setting aside the random walks L and R, the overlaps actually *decrease* with added mechanisms! As with the successor and ghost visit coverage metrics, there appears to be a particularly sharp drop with configuration 4, when agent preferences become active. Again, the power of HGNs in drawing agents together is clear in the increased overlap in configuration 6, but faces interference from influence edges in configuration 7.

The overall negative correlation between successor overlap and number of mechanisms is surprising. Perhaps the mechanisms lead the agents into parts of the CEG that they otherwise would not visit. Preferences in particular can lead agents to prefer highly branching regions that otherwise would be relatively inaccessible. In such a region with high node degree, SCAMP's stochastic roulette selection can push different runs in different directions, increasing successor coverage and thus reducing overlap. Modelers who assign favorable features to some events may unconsciously focus more attention on them and ramify the paths to which they lead more than they do for other events, a form of modeling bias of which they should be aware.

## 5.4   Impact of Individual Mechanisms

Table 1 compares the effect of preference and HGN mechanisms, aggregated over all configurations, against the baseline. In this table, an entry of the form " >, 0.04" means that the row variable is larger without the column mechanism than with it, with significance $p = 0.04$. "<" means that the unconstrained value is smaller than the constrained one, and "NS" means that $p > 0.05$. The aggregate impact of influence edges for all variables is NS.

Preference has the most widespread impact of all mechanisms, affecting every measure except ghost overlaps, and it reduces all measures that it affects except overlap among avatar visits, which it increases. The increase of avatar visit overlap is in line with our initial hypothesis: the more mechanisms constrain the model, the fewer different nodes agents will visit, and the more those sets of nodes will resemble each other across runs. The decrease of successor overlap is puzzling, but confirms the impression we drew from Fig. 7.

**Table 1** Mechanism impact on coverage

| Variable | Preferences | HGNs |
|---|---|---|
| Ghost visits | >, 7E−7 | >, 1E−5 |
| Ghost overlaps | NS | <, 0.01 |
| Avatar visits | >, 2E−4 | NS |
| Avatar overlaps | <, 1E−9 | <, 1E−9 |
| Successors | >, 7E−7 | >, 2E−5 |
| Successor overlaps | >, 8E−8 | NS |

HGNs are the next most influential mechanism, increasing visits (except for avatar visits) and decreasing overlaps (except for successor overlaps).

Though in the aggregate influence edges do not have significant impact on these measures, they can reduce the contribution of HGNs by limiting agent access to types of events that the HGN identifies as urgent.

# 6 Discussion and Future Work

Our results, though based on experiments with SCAMP, are important for the responsible use of any agent-based modeling framework, in two ways. (1) Modelers have a sense of the range of possibilities covered by their models, based on their static structure. The model's actual attractor when it runs may be much smaller. Users need to understand the effective coverage of a model under different conditions, and modelers need to understand how adding mechanisms is likely to impact that coverage. Sometimes users will want to increase coverage to consider more possible outcomes; in other cases they will want to decrease it to focus on the most likely outcomes. (2) Adding mechanisms to capture more dimensions of the real world can not only provide a more focused result, but also shift the location of that result in state space.

These results suggest several directions for future work.

- Validate our hypotheses about what features of SCAMP (e.g., widely varying branching factors?) lead to premature focusing, and develop guidelines for modelers who use SCAMP to avoid unrealistic expectations about how aware the system actually is of all the alternatives they are constructing.
- Explore in more detail what leads to some of the counterintuitive behaviors we have discovered, such as the interaction of HGNs and influence edges, and the decrease in successor overlap as we add mechanisms.
- Formally, a system's attractor is the region of state space to which it is constrained *after initial transients have died out*. Our data comes from complete runs, and ignores possible noise from start-up conditions. The start-up period can be identified by plotting the entropy of the roulette constructed by each agent as a function of time [12]. Applying this measure to the analysis in this paper is more challenging than in our previous application of it, but would refine our results.
- The notion of an attractor is only one of several physics-based concepts that can elucidate the dynamical behavior of a social simulation. We are exploring others, such as the graph spectra of emergent social networks.
- At several points, coverage and overlap measured by avatar visits behave very differently than ghost visits and successor counts. Most reports we generate for users concern the movement of avatars, and we have viewed the ghost mechanism

and successor structure of the CEG as internal details that are not relevant to analysts, but clearly they are important in assessing the model's dynamic coverage, and we will explore ways to communicate this information to users.

# References[2]

1. Brueckner, S., Parunak, H.V.D.: Information-driven phase changes in multi-agent coordination. In: Proceedings of Workshop on Engineering Self-Organizing Systems (ESOA, at AAMAS 2005), pp. 104–119, Springer (2005)
2. Butner, J.E., Wiltshire, T.J., Munion, A.K.: Modeling multi-agent self-organization through the lens of higher order attractor dynamics. Front. Psychol. **8**, 380 (2017)
3. Cenek, M., Dahl, S.K.: Geometry of behavioral spaces: a computational approach to analysis and understanding of agent based models and agent behaviors. Chaos Interdisc. J. Nonlin. Sci. **26**(11) (2016)
4. Falandays, J.B., Smaldino, P.: The emergence of cultural attractors: An agent-based model of collective perceptual alignment. In: Annual Meeting of the Cognitive Science Society, vol. 43 (2021)
5. Gilbert, N., Troitzsch, K.G.: Simulation for the Social Scientist, 2nd edn. United Kingdom, Open University Press, Buckingham (2005)
6. Huberman, B.A., Glance, N.S.: Evolutionary games and computer simulations. Proc. Nat. Acad. Sci. USA **90**(16), 7716–7718 (1993)
7. Kahneman, D., Tversky, A.: The simulation heuristic. In: Kahneman, D., Slovic, P., Tversky, A. (eds.) Judgment Under Uncertainty: Heuristics and Biases, pp. 201–208. Cambridge University Press, Cambridge, UK (1982)
8. Mudigonda, S., Núñez-Corrales, S., Venkatachalapathy, R., Graham, J.: Scheduler Dependencies in Agent-Based Models: A Case-Study Using a Contagion Model. Computational Social Science Society of the Americas, Springer, Santa Fe, NM (2021)
9. Núñez-Corrales, S., Friesen, M., Srikanth, M., Venkatachalapathy, R., Graham, J.: In-Silico Models with Greater Fidelity to Social Processes: Towards ABM Platforms with Realistic Concurrency. Computational Social Science Society of the Americas, Springer, Santa Fe, NM (2020)
10. Parunak, H.V.D., Brueckner, S.: Concurrent modeling of alternative worlds with polyagents. In: The Seventh International Workshop on Multi-Agent-Based Simulation (MABS06, at AAMAS06), pp. 128–141, Springer, Hakodate, Japan (2006)
11. Parunak, H.V.D., Morell, J.A., Sappelsa, L., Greanya, J.: SCAMP User Manual. Parallax Advanced Research, Beavercreek, OH (2020). https://www.abcresearch.org/abc/papers/SCAMPUserManual.zip
12. Parunak, H.V.D.: Learning Actor Preferences by Evolution. Computational Social Science (CSS21), CSSSA, Santa Fe, NM (2021)
13. Parunak, H.V.D., Greanya, J., McCarthy, M., Morell, J.A., Nadella, S., Sappelsa, L.: SCAMP's Stigmergic Model of Social Conflict. Computational and Mathematical Organization Theory (2021)
14. Sappelsa, L., Parunak, H.V.D., Brueckner, S.: The generic narrative space model as an intelligence analysis tool. Am. Intell. J. **31**(2), 69–78 (2014)
15. Savit, R., Brueckner, S.A., Parunak, H.V.D., Sauter, J.: General Structure of Resource Allocation Games. Altarum, Ann Arbor, MI (2002). https://www.abcresearch.org/abc/papers/RAGpaper.pdf

---

[2] Publications by the author are available at https://www.abcresearch.org/abc/papers.

16. Shnerb, N.M., Louzoun, Y., Bettelheim, E., Solomon, S.: The importance of being discrete: life always wins on the surface. Proc. Natl. Acad. Sci. USA **97**(19), 10322–10324 (2000)
17. Tukey, J.W.: Exploratory Data Analysis. Addison-Wesley (1977)
18. Wilson, W.G.: Resolving discrepancies between deterministic population models and individual-based simulations. Am. Nat. **151**(2), 116–134 (1998)
19. Wolfram, S.: Cellular Automata and Complexity: Collected Papers. Reading, MA, Addison-Wesley (1994)
20. Zia, A., Koliba, C.: The emergence of attractors under multi-level institutional designs: agent-based modeling of intergovernmental decision making for funding transportation projects. AI & Soc. **30**(3), 315–331 (2015)

# All the Right Moves? Systematically Exploring the Effects of Random Movement in Agent-Based Models

**Edmund Chattoe-Brown**

**Abstract**  Movement in Agent-Based Models, particularly so called random movement, frequently seems to be treated as a black box. This situation implies that implementation details don't matter to model outcomes. This chapter demonstrates the lack of concern with details of random movement processes using a literature review as a case study, shows that unreported implementation details can actually matter dramatically to model outcomes and considers the wider implications of these two related findings both for future research and modelling practice. The effect of movement assumptions is explored using the "switchable" approach in which pairs of models differing only in the aspect of interest are directly compared. The chapter shows how some Agent-Based Models may fail to distinguish between movement processes which effectively ensure full population mixing and those which do not. In the case where there isn't full mixing, details of implementation are shown to matter considerably to model behaviour. Taken together, the literature review and model analysis show that there is significant room for a more systematic and evidence based analysis of the role that movement assumptions play in the reported outcomes of Agent-Based Models.

**Keywords**  Agent-Based model · Random movement · Research methodology · Wolf sheep predation model · NetLogo

## 1  Introduction

It is a commonplace of Agent-Based Modelling that the behaviour of a model can depend considerably on its assumptions. Subject to the value of the "Preferred Proportion" parameter, the basic Schelling model of type segregation [1, 2] can display "loose" clusters with different agent types directly adjacent, "tight" clusters with different types needing to be separated by empty sites or a "rolling boil" of non-convergence [3] all but indistinguishable from the initial distribution. Exploration of

---

E. Chattoe-Brown (✉)

School of Media, Communication and Sociology, University of Leicester, Leicester, UK
e-mail: ecb18@le.ac.uk

553

the effects of parameter variation is often known as *sensitivity analysis* [4]. Similarly, but unfortunately with less agreement on terminology, we recognise that an Agent-Based Model with random mixing, for example, may display different properties than one with interaction through social networks, whether static or dynamic [5]. This might be considered the sensitivity of a model to *specification*—what elements the model includes—rather than to the values of parameters taking the model specification as given. So once you have decided to include a social network in your model, you need a parameter for the number of ties that each agent has but if you assume random mixing you don't.

This chapter considers a particular aspect of specification, namely the nature of agent *movement* (and specifically so called *random* movement) in Agent-Based Models. Examining a sample of published research as a case study suggests that (as widely feared but not always explicitly demonstrated) the implementation detail of such movement is often hard to access and a key result of this chapter is to show that, since the specific implementation can matter considerably to model outcomes, it needs to be properly documented. Standards are thus important not just for their own sake but because their absence can be shown potentially to undermine model credibility. This specific analysis therefore serves as a microcosm for the more general issue of adequately documenting Agent-Based Models so they—and their results—can be relied on for future use. But, of course, any conclusions drawn here are subject to qualification by subsequent literature reviews or model analyses that have a different basis or are more extensive. (For this reason, it is not appropriate to object that I am making claims about all models or all kinds of movement. I only present an evidenced hypothesis whose limitations are clearly stated. However, research has to start somewhere and I have not yet discovered any similar analyses in the literature to that presented here.)

The first section of the chapter reports a case study literature review of articles in the *Journal of Artificial Societies and Social Simulation*—hereafter JASSS—which mention random movement. This review is used to present and clarify the concepts that will be relevant for subsequent analysis and to support the claim that treating random movement as a black box is a significant tendency. The second section uses a standard NetLogo Agent-Based Model—the Wolf Sheep Predation model [6]—to show how the exact specification of so called random movement—which actually includes a range of relevant possibilities—can have a big impact on model outcomes. This illustrates why research in which movement is a black box may suffer from compromised credibility. The final section discusses these results and draws conclusions from the two related strands of argument. These include some recommendations for improved modelling practice based on evidence from the literature review and model analysis.

## 2 Literature Review

I use the term *enacted practice* to refer to what is actually done in published research as opposed to what is claimed in general or methodological discussion. Famously, Agent-Based Modelling presents itself as a potentially empirical method but its enacted practice is still largely to build models without data [7, 8]. My intention here is to report a case study of the enacted practice in a community of modellers regarding random movement, that is to distinguish what they actually do in publishing their research from what may be said about it in principle or in general.

The review was carried out by searching JASSS for two specific terms—"random movement" and "move randomly"—and reading all the articles (27) which included either of these terms. No article included both terms. There is insufficient room to describe these articles in the current chapter but a list of their relevant features—such as why the model code usually cannot be practically accessed—is available on request from the author. The aim of this approach is not to claim that this sample—which is actually a population for the designated search terms in JASSS—represents the whole field of Agent-Based Modelling or that there could not be models involving random movement *not* discovered by these search terms or that there is no research which does *not* treat movement as a black box. Instead, it is to support claims about a strong general tendency in enacted practice initially using a sample that is manageable for analysis and reporting in a single chapter along with model experimentation. Furthermore, it is not implausible that JASSS, as a long running journal focusing heavily on Agent-Based Models, should broadly manifest larger trends in such research generally. Be that it may, however, the main aim of the literature review is to provide an inductively grounded basis for a firmer conceptual understanding of random movement and a better justified statement of relevant issues for future research. Having stated a clear position based on evidence, future research can then confirm or refute my finding using larger or different samples [9, 10]. What is important for progressive science is not that each piece of published research be the last word—because that might be impossible or at least practically very difficult in terms of resources like research time and word count—but that it is clear and systematic enough for future research to develop.

The first thing to note is that splitting the time period during which JASSS has been published into two roughly equal parts—1998–2009 and 2010–2022 with the final year obviously being incomplete—does not suggest that random movement models selected by these particular search terms are becoming significantly less common. There are 13 instances in the first period and 14 in the second. This provides at least some justification for the approach of the present chapter. Random movement models would be less worth investigating if they were plainly dying out.

Secondly, we can make one major distinction and some minor ones from the articles reviewed. The major distinction is that at least two things can reasonably be meant by random movement. One occurs in *cellular* models, those which are arranged on a grid, where agents can only move to empty grid sites. Here randomness means, if any conditions for movement are satisfied, to identify all possible empty sites and

choose one at random with an equal probability. Even here, however, exact speci-
fication detail is important. The standard NetLogo implementation of the Schelling
model [11] actually *doesn't* do this but instead allows agents to keep moving *until*
they find an empty site. This potentially biases the agents towards shorter relocation
distances. The question is then, apart from issues of empirical realism, whether the
exact implementation matters to the reported model outcomes. Some evidence is
provided later in this chapter that it might but I have never seen the implications of
these variant assumptions—how the model is typically described and how it actually
runs—directly compared for the Schelling model. The other use of random move-
ment occurs in *free space* models where agents literally do wander around on a
geographical terrain. Here, unlike the cellular case, fully random movement is not
uniquely defined: It could involve (as we shall see) different movement distances,
different changes of orientation, different kinds of movement heterogeneity and so
on. Other more minor distinctions between reviewed models are also useful, for
example between models with general random movement and those only displaying
it under certain conditions, between models used directly to describe the social world
and those for other purposes—for example to develop modelling tools, discuss or
replicate the work of others and so on—and between worlds with "stuff" in them,
like food which agents forage, and worlds consisting only of other agents. In a way,
from this perspective, the most surprising thing about the basic Schelling model is
that it includes no real geography and that this fact is not widely remarked upon
[12]. The models containing stuff are often associated with deliberate movement in
at least some circumstances—if you can't see any food move at random, if you can
move towards it—while the models solely involving agents more commonly are not.

Thirdly, a focus on enacted practice in the specific area of movement illustrates that
general concerns about the reproducibility of Agent-Based Models are well founded
[13]. The issue is that, whatever the stated ideal, the great majority of the models
reviewed cannot, in fact, be traced back to their movement assumptions without
significant effort: The movement algorithm is not explained or only sketched out
roughly in pseudo code, actual access to the code is not mentioned at all, the code
is written in an outdated or obscure language, the code is only available from the
author—who may not now be contactable, the article claims to provide code but the
link is broken, the code won't run on current versions of NetLogo—so one has to trawl
back to the last version they will run on, the code won't run without undocumented
libraries or subsidiary files and so on. Some of these hands-on issues suggest concrete
improvements to practice. For example, should the file names for NetLogo code
always include a version number? If, as I shall show, the details of movement can
matter significantly to model outcomes, then not being able to reconstruct what was
actually assumed makes the results of most models in my review set potentially
non robust. More generally, the documentation and preservation of code needs to be
sufficient not just for present concerns but also for unanticipated future ones. There
is always a chance, as here, that we will need to revisit old models with new research
questions or objectives. Finally, as well as not being able to access the code, in many
(particularly older), models you can't see the actual simulated world in the article
but only graphs summarising run statistics. This may make it harder to get a sense

of what random movement actually involves. Compare the potential good practice in [14, Fig. 9] which shows exactly what the random movement of a single agent in the specific simulated world looks like. Note that I have done this analysis from the perspective of a single programming language, NetLogo, though this is widely used in JASSS and in the articles in my literature review. Nonetheless, many of the same issues about accessibility, like broken links or old versions, apply to other languages that I am not competent to evaluate. But my claim that a language is obscure is not merely subjective. It refers to how often it is used, whether it is still being maintained and updated and so on. I am not claiming that SWARM, for example, is obscure just because I don't know and use it personally.

Thus, as it is possible to trace very few models through to their actual movement assumptions, so it is particularly useful to discuss one of the few exceptions [15] in my sample in more detail. In this article, the movement procedure is "rt random 50 lt random 50 fd 1". Interestingly, this is exactly the same code as the Wolf Sheep Predation model whose properties I shall analyse further later in the chapter, showing that there is a danger that code may be reused incautiously without the possible implications being fully understood. This example provides a useful starting point for focused discussion of variations in possible movement assumptions which might reasonably be considered random. The first two instructions are stochastic (because of the random procedure) and allow an agent only to change its orientation broadly within its "line of sight" since agents in NetLogo have a direction in which they face defined as an attribute. The combination of rt and lt procedures means that an agent can reorient to the right up to 49° from line of sight and then reorient left up to the same amount. In NetLogo, random 50 generates an integer from 0 to 49. Since the random procedure draws from a uniform distribution a reorientation 37° to the right is just as likely as a reorientation 12° to the right. The net effect of both instructions, each with its distribution of outcomes, is that small net changes of orientation relative to the existing direction of travel are most common. To change net orientation by 48° to the right requires *both* an atypically large value of the right shift *and* an atypically small value of the left shift. Another way to see this point is to recognise what changes in orientation this combination of procedures *excludes*: An agent cannot simply turn on its heel, that is reorient 180°, in a single movement for example. The third procedure says that, after reorienting itself, the turtle moves forward a fixed amount which is only a fairly small distance relative to the world size. This example allows us to add precision to the distinction between calibration and specification discussed above. If forward movement is given as an integer, as it is here, then it can be considered as a parameter for calibration. But whether or not an agent reorients (and thus whether parameters are needed for the extent of that reorientation) is clearly a matter of specification. There is no reason in principle why agents should not simply move in straight lines with different initial orientations though I have never seen an assumption of that kind used. Whether a model which allows all possible reorientations should be considered differently specified to one that only allows some and favours small changes in orientation is a slightly less clear cut case. Finally, the reason why these variations should all be considered as examples of random movement is that the distribution of end positions relative to

start positions is clearly randomised. Using the procedure fd 1 will obviously expose an agent to different end state possibilities when compared to the procedure fd 2 but, in both cases, these final positions are not deterministic relative to the starting position despite the fixed integer in the fd command.

Based on analysis of existing research, I investigate the dependence of model outcomes on specific random movement assumptions. I have shown that movement as a black box is a distinct tendency in published research (admittedly from a single journal albeit an important one). Now I will show why that might be a significant problem.

## 3  "Switchable" Models of Random Movement

Like many developments in methodology, an idea can be extremely simple—almost ponderously so—as long as it is the right one. So, if we want to explore the effects of specification assumptions we simply build variant models differing only in these assumptions and compare them [5]. Given this approach, any differences observed in output behaviour *must* be the result of specification effects because there is no other possible source for them. This is the so called "switchable" procedure, analogous to sensitivity analysis for parameters but involving discrete rather than potentially continuous variation in model alternatives compared.

I will now demonstrate such an analysis on the Wolf Sheep Predation model [6] to support my earlier claims about the problems that could arise when random movement assumptions cannot be accessed and/or are not examined for their effects: A situation which applies in nearly all the articles I examined in my literature review as discussed above. This model was partly selected for analysis because it is included in the permanent NetLogo models library so my results can easily be checked and improved on. The exact version of the code used in this chapter is also available from the author to aid reproducibility. The basic model in the library also contains the specific movement procedure "rt random 50 lt random 50 fd 1" found in [15]. Apart from its simplicity and accessibility, this model was also chosen because it is fundamentally based on movement, it doesn't converge in a banal way—as a surprising number of Agent-Based Models seem to do given their claims about complexity—and the movement process applies to both sheep and wolves, thus maximising the possibility for interesting variations in emergent behaviour. Note, however, that this model was *not* chosen to be realistic, either as a description of sheep interacting with wolves or to stand for other Agent-Based Models aspiring to greater realism. The point here is solely to display significant specification variability from assumptions that may appear innocuous and thus not be described or analysed in typical research. It will necessarily be a topic for future research whether this variability proves to be widespread in diverse models including those which aspire to be realistic. But at present very few models of any kind appear to be tested for the sort of specification variability I demonstrate here.

To facilitate more runs of each simulation condition and thus better characterise the stochasticity of the model. I characterised each run in terms of *sampled* wolf and sheep populations. By sampled I mean that in each tick, I recorded these populations with a 10% probability and then, at the end of the run I summed this sampling list and divided by its length. This approach is intended to characterise each simulation run as a whole—rather than, say, just the average for the last 200 ticks—and thus accommodates the known situation with this model that individual runs may spend different lengths of time doing Lotka-Volterra cycles [16], displaying run in behaviour and so on while nonetheless displaying quite strong qualitative similarity overall. This means that although cycling may occur in different parts of the run, all runs contain significant cycling. The logic is that, as long as the sample of populations is not too small, it can give an overall sense of how two simulation conditions differ when coupled with repeat runs within a condition to assess stochastic variability. I leave the investigation of different—perhaps improved—characterisations of this model to future researchers as here I just want to clearly illustrate the main point about the dependency of model outcomes on the specific implementation of movement. Each experimental condition was run ten times and the typical maximum and minimum population values are within 7% of the mean suggesting that this system doesn't have problematic distributional properties requiring more repetitions than this in each condition. The logic of the experiments presented here is that if the mean population differences *between* conditions are much larger than the stochastic variation *within* conditions then it is reasonable to postulate that the difference is a meaningful effect of changing the movement assumptions.

The first result is to investigate the effect of different deterministic fd commands while keeping everything else fixed: Average wolf and sheep populations are reported to the nearest unit. Note that, because of randomness in reorientation, these runs are still stochastic and create random movement even though fd is fixed and deterministic.

As Table 1 shows, fd 1 behaves clearly differently from fd 5 and fd 10 while fd 5 and fd 10 are not plausibly distinguishable given the stochastic variation in individual run outcomes of about 7% for a single condition. Reflection suggests why this might be in that random movement is *not* random mixing. Random mixing involves an equal chance of each agent encountering all others but small fd values mean the chances of encountering distant agents are actually much lower. There may be a further complication here. In this model, it doesn't matter which sheep a wolf actually eats—distant or close. But in an epidemic model, for example, it *does* matter whether or not a susceptible agent is in contact with an infected one. It seems that bigger fd values "saturate" the interaction process: Going from fd 5 to fd 10 adds nothing appreciable to the outcome. Further, the fd value needed to saturate a population seems likely to depend on the size of the world, the number of agents and perhaps the size of patches: A patch in NetLogo being the unit of spatial extension which defines the size of the world. These values and their interpretation are not typically given a lot of attention in published models but the results in this chapter show that perhaps they need to be. Again, the style of analysis shown here needs to be generalised to other classes of models before its hypotheses can be taken as confirmed. In an epidemic model one can simply increase the density of agents by

**Table 1** Wolf and sheep populations with different forward movement distance

| Forward movement distance | Average wolf population | Average sheep population |
|---|---|---|
| 1 | 73 | 163 |
| 5 | 99 | 141 |
| 10 | 100 | 144 |

**Table 2** Wolf and sheep populations with "line of sight" and "completely random" reorientation

| Orientation assumption | Average wolf population | Average sheep population |
|---|---|---|
| rt random 50 lt random 50 | 73 | 163 |
| Random 360 | 3 | 205 |

assumption. In the Wolf Sheep Predation model, the sustainable population of wolves is dependent on the sustainable population of sheep and this, in turn, is dependent on the rate at which grass regrows. Thus population densities are not controlled directly by corresponding parameters. Put another way, certain wolf and sheep densities may not actually be achievable simply by changing the grass regrowth rate.

This hypothesis about saturated and unsaturated systems is given strong support by the next experiment. In the Wolf Sheep Predation model found in the NetLogo models library it is assumed that the direction of travel favours line of sight as discussed above. It might be argued that the procedure rt random 360 was closer in spirit to "true" random movement than rt random 50 lt random 50. Table 2 shows the effect of these two different reorientation specifications on wolf and sheep populations with everything else in the model remaining exactly the same.

Thus, in the region of *unsaturated* movement assuming fd 1, a change in the orientation assumption totally changes the model outcome basically flat lining the wolf population. This dramatic sensitivity to specification either requires a modeller to be sure their model is in the region of saturated movement—where further details of implementation apparently become largely irrelevant—or to justify the specifics of their assumptions in an unsaturated model so their results can be relied on. Interestingly, only one article in the set I reviewed seemed to provide any proximate justification for the random movement assumption and that was basically that the authors "felt" it was an adequate approximation for the domain. No reasons or evidence were given.

This view of saturated and unsaturated rates of movement is confirmed by introducing the simplest form of heterogeneity into the model. Instead of a constant fd value, this now becomes a *range* from a uniform distribution: For example, wolves either move forward 1 *or* 2 units with equal probability. As Table 3 shows, when the average of the movement range is in the saturated region, heterogeneity has almost no impact as before but when it is in the unsaturated range it makes at least some difference—though the scale of the effect is of the magnitude of the stochastic variability in runs. In fact, this small effect is an artefact of the unlucky circumstance

that saturation occurs in the Wolf Sheep Predation model with almost any other fd value *except* the one chosen as the default in the published code! This can easily be confirmed by looking at the behaviour of rt random 50 lt random 50 fd *2* which gives an average wolf population of 100 and sheep population of 149 which virtually indistinguishable, allowing for stochastic variation, from the fd 10 outcome. Thus the effect of heterogeneity is not intrinsically small but it is hard to display with this particular model and associated parameter values because there isn't room for move ranges that doesn't intrude into the saturated region.

This view can be confirmed (for use in future analysis) by lowering the rate at which grass regrows (the previous value of 30 is changed to 45). This lowers the sustainable sheep population and (in turn) the corresponding sustainable wolf population.

Table 4 shows that movement is not saturated at fd 2 as it was with the previous grass regrowth rate but only at fd 5 which is indistinguishable, given the stochastic variability in runs from fd 4. Interestingly, however, this analysis contradicts the earlier presumption that model behaviour is divided only into saturated and unsaturated zones. With fd 10 the populations resemble those for fd 3 suggesting it may be possible to move too *much* for random mixing as well as too little. However, analysis of this slightly unlikely looking result will have to be postponed.

In the next section I turn to some wider implications of this sort of analysis.

**Table 3** Wolf and sheep populations with different heterogeneous movement ranges

| Forward movement range (uniform distribution) | Average wolf population | Average sheep population |
| --- | --- | --- |
| 1–2 | 96 | 148 |
| 1–12 | 100 | 136 |

**Table 4** Wolf and sheep populations with different movement rates under a lower rate of grass regrowth

| Forward movement distance | Average wolf population | Average sheep population |
| --- | --- | --- |
| 1 | 24 | 156 |
| 2 | 45 | 143 |
| 3 | 50 | 139 |
| 4 | 61 | 130 |
| 5 | 62 | 129 |
| 10 | 52 | 141 |

## 4   Discussion and Conclusions

This chapter is a work in progress. Nonetheless, two main results have clearly been demonstrated worthy of further study. The first is that the exact specification of random movement can have a major effect on model outcomes. The second is that, following from this, we cannot allow a situation where the actual movement assumptions of many models—as evidenced by the literature review—can only be accessed with difficulty. I have proved that it is simply not the case that movement can safely be left in the black box where it predominantly seems to reside.

As a deliberately provisional—but clear and empirically supported—hypothesis, it seems to be the case that, in the saturated region, specific implementation details of random movement, whether fixed or heterogeneous, line of sight or completely random orientation, have no convincing effect on outcomes at least for the Wolf Sheep Predation model. Thus one practical strategy to increase model robustness is simply to report different fd values to demonstrate that the movement effect is saturated. If, for any reason, this strategy is not considered suitable, then further analysis should be done on the effect sizes of different movement assumptions and this creates a need for empirical calibration. To take the Wolf Sheep Predation model as an example, how fast do sheep and wolves typically move over a landscape in *real* units?

In fact, there are more variant assumptions that could still reasonably be considered random movement and whose comparative effects could usefully be investigated. Firstly, rather than draws from movement ranges the same for all breeds, agents could be heterogeneous in attribute terms so that some sheep and wolves are always faster, or faster on average, than others. A breed is simply a structure in NetLogo that allows agents to both share attributes—for example energy levels from food—but also have distinctive ones—so wolves might share meat while sheep could not share grass. Secondly, models could also be heterogeneous in breeds so that wolves and sheep do not draw from the same movement distributions. It is possible to avoid conflating the effects of more movement in the system generally by partitioning a fixed supply of movement so the outcomes are analysed from fd 9 for wolves and fd 1 for sheep all the way to fd 1 for wolves and fd 9 for sheep. The same can be done in comparing distributions of possible movement values with a fixed movement value that has the same mean. Thirdly movement could involve a chance in each tick that an agent stays still. This would correspond to a sheep moving to find fresh grass and then stopping to eat it. This last assumption might be more realistic than constant movement in some settings. This in turn suggests other movement processes that might be more plausible for some specific cases—like movement with regular returns to a base location—while still being reasonably described as random.

The point here is not to exhaustively characterise all forms of movement but merely to make clear—since it has been shown that it may make a large difference to model outcomes—that random movement doesn't unambiguously define a single plausible code implementation in all possible domains. In addition, and for similar reasons, it will probably prove necessary to explore the joint effect of world size, patch size and

population not only on the outcomes of the model but on the size of the saturated movement region. My next aim is to use the larger unsaturated region I created in the Wolf Sheep Predation model—by lowering the grass regrowth rate—to make a more effective exploration of the effect that different types of heterogeneity have on model outcomes and also to look more carefully not just at average populations but at the actual dynamics of their change i.e. whether different qualitative—or even quantitative—regimes can be discovered in the changes of wolf and sheep populations.

Thus the enacted and inductive approach to existing models proposed in the literature review not only identifies real rather than conjectured problems to fix but also suggests operational rather than abstract approaches to investigation. If you want to know what matters in a particular model, don't guess but investigate using the switchable approach and report your results so they can be challenged/extended. If you want to know how to improve model reporting procedures, look at exactly what goes wrong when you try to run down actual code from a sample of models.

This approach also suggests other forms of analysis that may be beneficial, specifically more investigation of single agent and paired agent trajectories and their properties. In the higher grass regrowth case, *why* does the fd 2 model saturate in this specific case while the fd 1 model does not? Why does line of sight movement potentially have such a different effect from truly random movement in a model that is otherwise identical? I suspect the behaviour of moving populations may be another case where simple intuition will fail us without the support of systematic model based analysis. While imagining the effects for two agents, we can too easily disregard the effects of all the others which will jointly be moving too. Could other things be happening in the typical black box of agent movement like mere bugs or, more worryingly, random number artefacts [17]. Another passing argument for removing movement from the black box is that it makes it easier to see whether what is happening in the model is actually what the modeller wants or expects. This may suggest that—like [14]—we should demonstrate the operation of random movement straightforwardly even if just in an appendix.

It should also be noted that this analysis involves a single programming language: NetLogo. It seems unlikely that implementations of random movement in SWARM or MASON would raise completely different issues but certainty requires investigation.

Finally, it is important not to miss the point of this chapter. The findings here do not become irrelevant simply because one takes the view that random movement is implausible or is in declining use in Agent-Based Modelling though a study of enacted practice suggests that this belief isn't actually supported by evidence. Exactly the same concerns can be identified for *any* movement procedure that remains in a black box. If agents move purposively, how much does it matter to the outcome if they jump to the new patch or travel there in finite time? Clearly it will matter, for example, in the case of epidemic spread. Does the step distance of non random movement or patch size still affect model outcomes as it has been shown to do with random movement? The definite conclusion I reached from reading the set of articles in the literature review was that the authors appeared to believe that the details of the random movement assumption didn't matter and/or weren't a problem. I have

shown that this belief may be imprudent and how we should go about investigating the effects of movement further in a wider class of models. This is the start of an investigation. It could not sensibly be its end.

This approach also has a wider implication for our attitude to progressive research. I know the patterns I have identified are only in JASSS articles and that they may not therefore generalise to all Agent-Based Models. It is also possible that I could have used better search terms to find more—or more relevant—articles featuring random movement. But these points do not actually—perhaps counterintuitively— detract from my results. Given the stated limitations of my research, my hypotheses and evidence stand and are clearly stated so they can be evaluated and developed. They may, of course, be undermined or qualified by subsequent research on different samples but we will not know until that is actually done. Only publishing research which can be the last possible word may actually be a barrier to scientific progress. Because in this case it would involve finding and reading a potentially very large number of articles which there might then not be space to report effectively. This could therefore be seen as an example of the *Nirvana fallacy*—comparing an actual piece of research with its ideal—but actually unrealistic—alternative rather than with another actual piece of research. The extent to which movement is a black box in the whole population of Agent-Based Models is not yet known, and might unfeasible to determine, but in this small but non trivial sample, the tendency is clearly considerable.

# References

1. Schelling, T.: Models of segregation. Am. Econ. Rev. **59**(2), 488–493 (1969)
2. Schelling, T.: Dynamic models of segregation. J. Math. Sociol. **1**(2), 143–186 (1971)
3. Chattoe-Brown, E.: Why sociology should use agent based modelling. Sociol. Res. Online **18**(3) (2013)
4. Ten Broeke, G., Van Voorn, G, Ligtenberg, A.: Which sensitivity analysis method should I use for my agent-based model? J. Artif. Societies Soc. Simul. **19**(1) (2016)
5. Chattoe-Brown, E.: Why questions like "do networks matter?" matter to methodology: how agent-based modelling makes it possible to answer them. Int. J. Soc. Res. Methodol. **24**(4), 429–442 (2021)
6. Wilensky, U.: Netlogo wolf sheep predation model. Center for Connected Learning and Computer-Based Modeling, Northwestern University, Evanston, IL (1997). http://ccl.northw estern.edu/netlogo/models/WolfSheepPredation. Last accessed 2022/06/03
7. Angus, S., Hassani-Mahmooei, B.: "Anarchy" reigns: a quantitative analysis of agent-based modelling publication practices in JASSS, 2001–2012. J. Artif. Societies Soc. Simul. **18**(4) (2015)
8. Dutton, J., Starbuck, W.: Computer simulation models of human behavior: a history of an intellectual technology. IEEE Trans. Syst. Man Cybern. **SMC-1**(2), 128–171 (1971)

9. Keijzer, M.: If you want to be cited, calibrate your agent-based model: reply to Chattoe-Brown. Rev. Artif. Societies Soc. Simul. (2022). https://rofasss.org/2022/03/09/Keijzer-reply-to-Chattoe-Brown. Last accessed 2022/06/03
10. Chattoe-Brown, E.: If you want to be cited, don't validate your agent-based model: a tentative hypothesis badly in need of refutation. Rev. Artif. Societies Soc. Simul. (2022). https://rofasss.org/2022/02/01/citing-od-models. Last accessed 2022/06/03
11. Wilensky, U.: Netlogo segregation model. Center for Connected Learning and Computer-Based Modeling, Northwestern University, Evanston, IL. http://ccl.northwestern.edu/netlogo/models/Segregation. Last accessed 2022/06/03
12. Chattoe-Brown, E.: Un drôle de type: the schelling model, calibration, specification, validation and using relevant data. In: Ahrweiler, P., Neumann, M. (eds.) Advances in Social Simulation: ESSA 2019, pp. 243–255. Springer, Cham (2021)
13. Donkin, E., Dennis, P., Ustalakov, A., Warren, J., Clare, A.: Replicating complex agent based models, a formidable task. Environ. Model. Softw. **92**, 142–151 (2017)
14. Costopoulos, A.: Evaluating the impact of increasing memory on agent behaviour: adaptive patterns in an agent based simulation of subsistence. J. Artif. Societies Soc. Simul. **4**(4) (2001)
15. Bao, H., Wu, X., Wang, H., Li, Q., Peng, Y., Lu, S.: Conflicts induced by different responses to land expropriation among the farmers involved during urbanization in China. J. Artif. Societies Soc. Simul. **22**(1) (2019)
16. Thierry, H., Sheeren, D., Marilleau, N., Corson, N., Amalric, M., Monteil, C.: From the Lotka-Volterra model to a spatialised population-driven individual-based model. Ecol. Model. **306**, 287–293 (2015)
17. Polhill, J., Izquierdo, L., Gotts, N.: The ghost in the model (and other effects of floating point arithmetic). J. Artif. Societies Soc. Simul. **8**(1) (2005)

# A Template for Transfer of NetLogo Models to High-Performance Computing Environments for Enhanced Real-World Decision-Support

**Jason Thompson, Haifeng Zhao, Sachith Seneviratne, Rohan Byrne, Rajith Vidanaarachchi, and Roderick McClure**

**Abstract** The sudden onset of the COVID-19 global health crisis and associated economic and social fall-out has highlighted the importance of speed in modeling emergency scenarios so that robust, reliable evidence can be placed in policy and decision-makers' hands as swiftly as possible. For computational social scientists who are building complex policy models but who lack ready access to high-performance computing facilities, such time-pressure can hinder effective engagement with end-users. Popular and accessible agent-based modeling platforms in computational social science such as NetLogo can make models fast to develop, but slow to run when exploring broad parameter spaces on individual workstations. However, while deployment on high-performance computing (HPC) clusters can achieve marked performance improvements, transferring models from workstations to HPC clusters can also be a technically challenging and time-consuming task for social scientists or those from non computer science-related backgrounds. In this paper we present a set of generic templates that can be used and adapted by NetLogo users who have access to HPC clusters but require additional support for deploying their models on such infrastructure. We show how model run-time speed improvements of between 200× and 400× over desktop machines are possible using (1) a benchmark 'wolf-sheep predation' model in addition to (2) an example drawn from our own applied policy modeling work surrounding COVID-19 management settings for Government in Australia. We describe how a focus on improving model speed is a non-trivial concern for model developers in the social sciences and discuss its practical importance for improved policy and decision-making in the real world. We provide all associated documentation in a linked git repository.

J. Thompson (✉) · H. Zhao · S. Seneviratne · R. Byrne · R. Vidanaarachchi
Transport, Health and Urban Design Research Laboratory, University of Melbourne (UoM), Victoria, Australia
e-mail: jason.thompson@unimelb.edu.au

R. Vidanaarachchi · R. McClure
Faculty of Health and Medicine, University of New England, Armidale, NSW, Australia

## 1  Background

In 2020, Australia's twin crises of the catastrophic 'Black Summer' bushfires and Covid-19 pandemic provided stark examples of crises that can be classified as X-events [1]; critical systems failures and crises that are at once extreme, sudden, novel, rare, surprising and disastrous. By definition, X-events have a relatively short unfolding time, but their impact is significant and may last decades or longer. As such, X-events hold important ramifications for the function of society and the intersection of science and policy.

A common feature of X-events is their association with the design and evolution of human, sociotechnical systems [2]. They may be exacerbated by human systems, emerge as a result of activity within human systems, or are crises of abstract, but critical human-designed systems that societies rely upon to function 'normally'. Such systems could include health care, banking and financial systems, communications, political, transport, economic and insurance systems.

As the world becomes increasingly connected through technological, social, geographic, and economic ties, the frequency of X-events is expected to accelerate [2, 3]. This is bad news. In Australia alone, the bushfire and pandemic-related crises [4] absorbed hundreds of billions of dollars in direct and indirect costs [5], plunged the nation into the worst economic crisis since the Great Depression, and destroyed environment, lives, livelihoods, and property at unprecedented scale [6]. In turn, the effects of both crises reverberated across numerous associated systems, exposing significant weaknesses and generating fresh crises within healthcare [7, 8], housing [9], education [10], transport [11], economics [12], finance [13], environment [6], politics [14], and industrial relations systems [15], each requiring their own adjustments and policy responses. Reasonably, many of these consequent effects were unforseen because no empirical record of their occurrence existed prior to the event.

In the early stages of the COVID-19 pandemic, the dearth of historical record and available data compromised scientists' typical means of gathering evidence about the world as well as their capacity to build models to deal with the crises and its corrolaries [16]. Because there was limited science to draw upon, it also compromised the application of evidence-based policy [17], which promotes 'following the science' [18–20].

In the absence of complete evidence, simulation modeling can provide a useful theoretical and practical bridge for scientists and policymakers, alike. For example, a branch of simulation—computational social science—is the discipline of representing communities, societies and social phenomena through the generation of tangible, observable, but computer-generated artificial or 'synthetic' societies. By authentically representing both known (e.g., from evidence) and proposed critical features, structures, and mechanisms of interaction among agents within artificial

societies, phenomena representing realistic potential crises befalling a society can be generated from the bottom-up [21]. Similarly, if crises within artificial societies can be generated, so too can policy solutions that prevent those same crises from arising. Simulation modeling of this type (primarily through agent-based modeling (ABM)) has demonstrated great utility across the world during the COVID-19 pandemic. It has proven to be a flexible, robust, and transparent tool that has provided valuable insight into the spatial, biological, and economic effects of the crisis and potential policy remedies [e.g., 22–30].

However, despite greater awareness of ABM and its strengths, challenges remain. Very simply, the current time it takes to develop, analyse, and iterate trusted models of artificial societies is often too long to make them useful to policy-makers. This delay can result in either (1) disengagement by time-poor policy-makers who require faster answers to 'what-if?' questions than is currently possible, or (2) the real-world crisis moves on to a new phase that is outside the scope of the current model. In both cases, 'the science' has failed to keep pace with decision-makers' needs [31].

There is therefore an urgent need for science and policy to connect better when faced with novel crises (e.g., COVID-19, environmental crises, and natural disasters) requiring up-to-date information and fast decision-making. Our own experience in working with policy-makers in both development and analysis of important social policy models demonstrates that the utilisation of HPC clusters is central to achieving this goal [20]. That is, once a set of policy-settings are agreed upon, the ability to run experiments, analyse, and feed-back results and insights quickly is critical.

In this paper, we demonstrate the advantages for modelers working at the intersection of computational social science and policy-making of deploying existing policy models developed in NetLogo [32] on HPC clusters featuring parallel computing infrastructure. Our aim is to assist computational social scientists, social science researchers and other regular users of the NetLogo software platform make a transition to using HPC clusters by providing a generic framework for adaptation by individual users through a set of step-by-step instructions and scripts. These can be used as-is or modified with the assistance of local expertise to suit researchers' own HPC environments.

## 2 Method

To demonstrate improvements in speed associated with the deployment of policy-models on HPC clusters, we used two examples. Firstly, the benchmark 'Wolf-sheep predation' model drawn from the standard NetLogo Models library. Secondly, we used a model developed in consultation with the Department of Health in Victoria, Australia to estimate risk associated with easing social restrictions after that state's 2nd wave of SARS-CoV-2 infections in 2020 [22, 29, 33, 34]. To demonstrate both compute and real-time performance differences between various HPC set-ups, we also compared run-times for the benchmark model when allocated to the HPC across 1, 2, 4, 8, 16 or 32 cores.

In example 1, we used the standard 'wolf-sheep predation' representation with minor adjustments so that it runs in NetLogo BehaviourSpace. Changes to the model include the addition of a global variable 'repetitions' on line 3 of the model code (see Sect. 2.1, below), as well as the removal of the text pop-up warning on line 59 that halts the model when it reached a 'max-sheep' threshold.

A BehaviourSpace function was then created called 'HPC_Experiment'. This function included 100 random numbers under 'repetitions', and also included 5 levels across each of the variables: 'wolf-gain-from-food', 'wolf-reproduce', 'sheep-gain-from-food', 'grass-regrowth-time', and 'sheep-reproduce' for a total of 312,500 individual model runs. The maximum time-step limit for each run was set to 150.

Example 2 compared a single scenario of 100 model runs under Policy 4 (aggressive elimination) from the authors' previously published and implemented COVID-19 epidemiological model [22, 29, 33, 34]. We ran the model for a total of 1500 time-steps, equivalent to 1500 model days.

Both examples were first run on one of the author's laptops (Intel® Core™, 4 cores, i7-7700HQ CPU @2.80 GHz, 32 GB RAM, Windows 10, 64bit OS). It was then deployed using the 'snowy' HPC partition on the University of Melbourne's 'Spartan' HPC cluster [35] using 8 cores per task. This is a traditional cluster with a high-speed interconnect in one partition as well as an alternative queue that uses virtual machines with a common image. Computing jobs are submitted to a Slurm workload manager specifying which partition they would like to operate on (e.g., in our case, 'Snowy'). Step-by-step information on how to prepare existing NetLogo models for deployment on the HPC, as well as the example benchmark model is contained at (https://github.com/melbhz/netlogo-hpc). A brief description of the procedure follows.

## 2.1 Description of the Procedure for Deploying NetLogo Models on the HPC

Firstly, NetLogo must be installed to run on the HPC cluster. Next a NetLogo model must be created to match the format required for deployment (described below and in the documentation). Any NetLogo dependencies and plug-in packages (e.g., rngs, GIS, etc.) should be copied to the same folder to the NetLogo model or can be placed in the 'extensions' folder in the NetLogo extensions directory.

Regardless of the type of model being run, it must contain a named experiment within NetLogo's BehaviourSpace function that contains a dummy input variable that bears no consequence to the function of the model (e.g., 'repetitions' from the example above). This dummy variable should contain a parameter space equal to the number of individual runs desired for each unique parameter combination. For example, the dummy variable 'repetitions' could contain a list of integers in the list [1:100]. Then, when combined with 3 policy setting choices on (for example) a real

variable #1 (low, med and high) and 5 choices on (for example) a real variable #2 (very low, low, medium, high, very high), this creates $100 \times 3 \times 5 = 1500$ individual model runs containing 100 runs of 15 separate policy combinations.
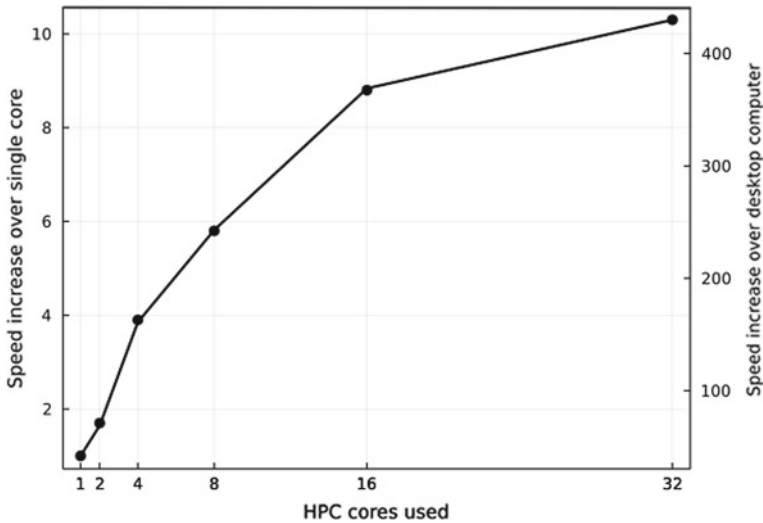
Using the set of templates available at (https://github.com/melbhz/netlogo-hpc), the user can open the 'create_xmls.sh' file, altering the highlighted inputs to match their own NetLogo file location, the unique BehaviourSpace model name, and their unique output file name location(s). Running this file in the command line will create a folder containing 100 separate.xml files that will be submitted as independent jobs across the HPC cluster and contain the full complement of parameter combinations. Computation speed in this case is ideally about 200 times faster (for example, if moving from 4 cores on a desktop machine to 8 cores on the HPC), however the actual computation time will depend on the number of CPU cores available to the experimenter on the HPC cluster at any time. As a practical example, during an academic teaching semester when students or other staff are also using your HPC cluster for their own work, you can expect a longer duration between deploying your model and having results returned because your job may be 'queued' behind others as you wait for computing resources to become available (see Results for further details).

Next, the experimenter can open and revise the 'submit_jobarray.slurm' file, again altering highlighted details related to file locations, input and output folders, as well as SBATCH settings unique to their own HPC environment. In our own case, we set: the number of nodes per job (we recommend 1), the desired partition on the HPC, the job name, the maximum run time for any job, and the desired CPUs per job (e.g., 8). In our specific case at the University of Melbourne, computing jobs that contain an additional 'critical' status are given priority over other jobs in the queue. Readers may also be able to request similar priority status from your local HPC administrators if applicable.

## 3 Results

Run time for our benchmark 'wolf-sheep predation' model showed that when deployed on the author's desktop computer, the complete set 312,500 trials was completed in 3 h and 35 min. When deployed on the HPC cluster using 8 cores, CPU time was instead 5 min and 57 s, while job wall-clock time was 53 s. Deploying on the HPC cluster resulted in a $243\times$ increase in real-time speed (See Fig. 1).

Performance improvements when deployed on the HPC were similarly improved for our real-world COVID-19 policy model. When run on the author's laptop, job completion took 2 h, 26 min. By comparison, CPU time for the same model deployed on the HPC using 8 cores was 26 min and 16 s, while job wall-clock time was 3 min, 17 s. Deploying the policy model on the HPC cluster resulted in a $44.5\times$ increase in real-time speed.

**Fig. 1** Comparative performance (speed) improvements gained over a single core allocation on the HPC cluster (left axis) and over the author's desktop computer (right axis), and for the benchmark 'wolf-sheep predation' model

Finally, comparison of performance improvements gained on the HPC cluster when the number of cores allocated to the benchmark task was manipulated indicated diminishing speed returns with increased allocated cores. As shown in Fig. 1, allocation of 8 cores (as used in the examples above) produced a $5.8\times$ speed increase over using a single core. Increasing to 16 cores increased performance by a further 51% to $8.8\times$ over the single core allocation. However, allocating 32 cores then only improved performance by an additional 17% to $10.3\times$ the speed of a single core (see Fig. 1).

It should also be noted that while wall-clock time was significantly improved in the 16 and 32 core conditions, in practice, allocation of the jobs to the HPC partition were delayed by 3 to 4 h because the busy HPC cluster needed to wait for sufficient allocation space to become available before deploying the model—this was despite being provided with priority access on the network. This result is non-trivial for real-world, time critical, policy modelling. It highlights (as described in Sect. 2.1) that it is very important to be mindful of the trade-off between time gained through efficiency of processing vs HPC cluster access when making decisions about how and when to deploy models on individual—and possibly congested—HPC infrastructure.

It is also important to consider that this experiment did not incorporate any run-time parallelisation load-balancing. The processing to be carried out was split at initialisation to run across the allocated number of cores in a static manner. A more dynamic allocation of cores at run-time would potentially improve the diminishing returns observed in Fig. 1. This is because during dynamic allocation the remaining processing is split across cores as soon as they finish the last workload assigned,

whereas in static allocation it is possible that a core that has finished its workload may idle until the end of the job while other cores finish theirs. In general, the decision for deploying on HPC infrastructure can be informed by such characteristic curves. The optimal setting for deployment needs to weigh the benefits of increased speed versus the costs and availability of cores on the HPC infrastructure.

## 4 Discussion

A hindrance to the uptake of computational social science and synthetic societies research and models to date has been the speed at which ABMs can be developed, explored, and run with timely results fed back to policy-makers. This was highlighted in the early stages of the COVID-19 pandemic where highly influential models in both the UK and Australia were adapted from existing influenza models rather than built as bespoke representations [36]. It shows that in urgent, unexpected crises where answers are demanded in minutes or hours rather than days (e.g., X-events), the capacity to ramp-up model speed and analysis is critical [20].

In our own work with the Victorian Department of Health during the second wave COVID-19 crisis in 2020 [28, 29, 34], extreme time pressure was exerted to match the timeframe of Victorian 'Crisis Cabinet' deliberations. Had our research group not been able to meet policy-makers' timelines, advice would likely have been omitted, resulting in arguably poorer decision-making, at least or decision-making that incorporated less information rather than more. After requested input parameter adjustments from Departmental officials covering wide parameter sweeps and the exploration of consequently large phase spaces, many model versions were run 'overnight' or over 24 h split across multiple individual and virtual machines before merging results. This delayed timely provision of advice back to Government and consequent decision-making. In turn, this affected millions of Victorians' lives as they waited for official Government advice on when and how social restrictions would be lifted in response to COVID-19 infection rates. It also provided motivation to build and share these HPC templates for adaptation and use by other researchers lest the value of insights provided by computational social scientists and other simulation modelers (e.g., epidemiological forecasters) be neglected on the basis of that advice being too slow to generate.

Returning to the concept of X-events, the rapid production and analysis of models facilitated through deployment on HPC clusters is not only potentially advantageous for faster, more informed decision-making, but also for the speed and evolution of models, themselves. Because results are available more quickly, HPC powered ABMs can be iterated and evolve with faster turn-around time and are more likely able to match real-time developments of X-events (e.g., natural disasters including bushfires and floods) in matter of minutes by incorporating new information and data as it comes to hand. This is more likely to enable ABM simulation to be relevant in the

face of rapidly unfolding crises that contain surprising or unexpected developments [1, 37]. Of course, exploring potentially more performant software platforms (e.g., Agents.jl) [38] and coding practices [39] should also be considered as a core strategy in this regard.

## 5   Conclusions

Simulation models used in supporting important public policy and decision-making should be robust [20, 40], but it is also important to recognise that sometimes decision-makers cannot wait for complete evidence before acting, especially in unfolding health crises or natural disasters [41]. The speed at which supporting synthetic evidence created through simulation modeling can be produced, analysed and presented is critical if computational social science is to keep pace with the world and deliver evidence in a form that directly addresses policy-makers' real-time requirements [31]. The sooner robust evidence can be presented, the sooner it has a chance to be incorporated into decision-making, and the greater chance it has to positively affect the course of strategy, policy direction, action and outcomes. In addition to other documented performance improvement measures that can be achieved in simulation modeling platforms [39], the utilisation of HPC clusters can assist to bring the production and presentation of important evidence generated by simulation modelers forward in time. Providing this capacity to social scientists by reducing barriers of access to HPC environments forms part of this effort.

## References

1. Casti, J.L.: X-Events: The Collapse of Everything. Harper Collins (2012)
2. Walsh, M.G., et al.: Whence the Next Pandemic? The Intersecting Global Geography of the Animal-Human Interface, Poor Health Systems and Air Transit Centrality Reveals Conduits for High-Impact Spillover. One Health, p. 100177 (2020)
3. de Ruiter, M.C., et al.: Why we can no longer ignore consecutive disasters. Earth's Future **8**(3), e2019EF001425 (2020)
4. Flannery, T.: The megafires and pandemic expose the lies that frustrate action on climate change. In: Fire, Flood, and Plague—Essays About 2020. The Guardian, Australia (2020)
5. Commonwealth of Australia: Budget 2020–21. Canberra, Australia (2020)
6. Binskin, M., Bennett, A., Macintosh, A.: Royal Commission into National Natural Disaster Arrangements. Royal Commission into Natural Disaster Arrangements, Australia (2020)
7. MacIntyre, C.R., Heslop, D.J.: Public health, health systems and palliation planning for COVID-19 on an exponential timeline. Med. J. Austr. **1** (2020)
8. Blecher, G., Blashki, G.A., Judkins, S.: Crisis as opportunity: how COVID-19 will reshape the Australian health system. Med. J. Austr. (2020.
9. Power, E.R., Rogers, D., Kadi, J.: Public housing and COVID-19: contestation, challenge and change. Int. J. Hous. Policy **20**(3), 313–319 (2020)
10. Drane, C., Vernon, L., O'Shea, S.: The Impact of 'Learning at Home' on the Educational Outcomes of Vulnerable Children in Australia During the COVID-19 Pandemic. National Centre for Student Equity in Higher Education, Curtin University (2020)

11. Beck, M.J., Hensher, D.A.: Insights into the impact of COVID-19 on household travel and activities in Australia—the early days of easing restrictions. Transp. Policy **99**, 95–119 (2020)
12. O'Sullivan, D., Rahamathulla, M., Pawar, M.: The impact and implications of COVID-19: an Australian perspective. Int. J. Commun. Soc. Dev. **2**(2), 134–151 (2020)
13. Andrew, J., et al.: Australia's COVID-19 public budgeting response: the straitjacket of neoliberalism. J. Publ. Budgeting Account. Fin. Manage. (2020)
14. Greer, S.L., et al.: The comparative politics of COVID-19: the need to understand government responses. Glob. Public Health **15**(9), 1413–1416 (2020)
15. van Barneveld, K., et al.: The COVID-19 pandemic: lessons on building more equal and sustainable societies. Econ. Labour Relations Rev. **31**(2), 133–157 (2020)
16. von Borzyskowski, I., et al.: Data Science and AI in the Age of COVID-19. The Alan Turing Institute: London, United Kingdom (2020)
17. Holmes, D., et al.: Deconstructing the evidence-based discourse in health sciences: truth, power and fascism. Int. J. Evid. Based Healthc. **4**(3), 180–186 (2006)
18. Mercuri, M.: Just Follow the Science: A Government Response to a Pandemic. Wiley Online Library (2020)
19. Abbasi, K.: Covid-19: politicisation, "corruption", and suppression of science. BMJ **371**, m4425 (2020)
20. Thompson, J., et al.: A framework for considering the utility of models when facing tough decisions in public health: a guideline for policy-makers. Health Res. Policy Syst. **20**(1), 107 (2022)
21. Epstein, J.M.: Generative Social Science: Studies in Agent-Based Computational Modeling. Princeton University Press (2006)
22. Blakely, T., et al.: The probability of the 6-week lockdown in Victoria (commencing 9 July 2020) achieving elimination of community transmission of SARS-CoV-2. Med. J. Austr. **213**(8), 349–351e1 (2020)
23. Blakely, T., et al.: Association of simulated COVID-19 policy responses for social restrictions and lockdowns with health-adjusted life-years and costs in Victoria, Australia. JAMA Health Forum **2**(7), e211749–e211749 (2021)
24. Kerr, C.C., et al.: Covasim: an agent-based model of COVID-19 dynamics and interventions. medRxiv, 2020.05.10.20097469 (2020)
25. Abeysuriya, R., et al.: Estimating risks associated with early reopening in Victoria (2021)
26. Chang, S.L., et al.: Modelling transmission and control of the COVID-19 pandemic in Australia. arXiv:2003.10218 (2020)
27. Rockett, R.J., et al.: Revealing COVID-19 transmission in Australia by SARS-CoV-2 genome sequencing and agent-based modeling. Nat. Med. **26**(9), 1398–1404 (2020)
28. State Government of Victoria: Emerging From Lockdown: Evidence, Modelling, Outputs and Assumptions. In: D.O.H.A.H. Services (ed.). State Government of Victoria, Melbourne, Victoria (2020)
29. State Government of Victoria: Emerging From Lockdown—Model. State Government of Victoria, Melbourne, Victoria (2020)
30. Milne, G.J., et al.: A modelling analysis of the effectiveness of second wave COVID-19 response strategies in Australia. Sci. Rep. **11**(1), 1–10 (2021)
31. Fischhoff, B.: Making decisions in a COVID-19 world. JAMA (2020)
32. Wilensky, U.: NetLogo Version 6.2.0. Centre for Connected Learning and Computer-Based Modeling, Northwestern University, United States of America (2021)
33. Thompson, J., et al.: The Estimated Likelihood of Eliminating the SARS-CoV-2 Pandemic in Australia and New Zealand Under Current Public Health Policy Settings: An Agent-Based-SEIR Modelling Approach. Available at SSRN 3588074 (2020)
34. Thompson, J., et al.: Modelling SARS-CoV-2 disease progression in Australia and New Zealand: an account of an agent-based approach to support public health decision-making. Austr. N. Z. J. Public Health **46**(3), 292–303 (2022)
35. Lafayette, L., Wiebelt, B.: Spartan and NEMO: two HPC-cloud hybrid implementations. In: 2017 IEEE 13th International Conference on e-Science (e-Science) (2017)

36. Squazzoni, F., et al.: Computational models that matter during a global pandemic outbreak: a call to action. J. Artif. Soc. Soc. Simul. **23**(2), 10 (2020)
37. Wilenius, M., Casti, J.: Seizing the X-events. The sixth K-wave and the shocks that may upend it. Technol. Forecast. Soc. Change **94**, 335–349 (2015)
38. Datseris, G., Vahdati, A.R., DuBois, T.C.: Agents.jl: a performant and feature-full agent based modelling software of minimal code complexity. arXiv:2101.10072 (2021)
39. Railsback, S.F., et al.: Improving execution speed of models implemented in NetLogo. J. Artif. Soc. Soc. Simul. **20**(1), 3 (2017)
40. Calder, M., et al.: Computational modelling for decision-making: where, why, what, who and how. R. Soc. Open Sci. **5**(6), 172096 (2018)
41. Thompson, J., McClure, R., de Silva, A.: A complex systems approach for understanding the effect of policy and management interventions on health system performance. In: Social-Behavioral Modeling for Complex Systems, pp. 809–831 (2019)

# Assessing the Cost of Population Dynamics Design Options in a Microsimulation

**Rachel J. Bacon** ⓘ **, George Hodulik** ⓘ **, and Wesley J. Wildman** ⓘ

**Abstract** We explore microsimulation design options as a source of divergence in total population when using demographic statistics from the United Nations to model population dynamics in three countries between 1950 and 2100. We compare 176 unique model designs, which toggle options such as the time step, the initial sample size of agents, variance reduction, ordering of demographic events, and adjustments to risk assignment as appropriate to each statistic. Results indicate that small population samples and 1-year time steps can produce particularly high divergence from UN targets, even when other options known to reduce divergence are implemented. Small sample 1-year models with low divergence are possible, but the specific combinations of options interact with a country's population dynamics in unpredictable ways, which prevents the design from being used in other country contexts. These findings are important for balancing efficiency, accuracy, realism, and generalizability in demographic microsimulation design.

**Keywords** Microsimulation · Demography · Model design

## 1 Demographic Statistics and Population Dynamics

### 1.1 Introduction

It is common for microsimulations, and some hybrid agent-based models, to use demographic statistics (e.g. fertility, mortality, and migration statistics) as exogenous variables to govern population dynamics in the model [1]. This is a straightforward way to maintain a realistic population size and composition, which strengthens the interpretation of model results as applicable to real-world populations [2]. Even

R. J. Bacon (✉) · G. Hodulik · W. J. Wildman
Center for Mind and Culture, Boston, MA 02215, USA
e-mail: rbacon@mindandculture.org

W. J. Wildman
Boston University, Boston, MA, USA

when using this seemingly simple approach, however, the model's births, deaths, and population can diverge from expected totals. This can happen if the demographic statistics are implemented incorrectly, but also for reasons that are not immediately obvious, such as a particular combination of mundane design options. We explore how selecting different microsimulation design options becomes a source of divergence when using demographic statistics from the United Nations (UN). Our analysis compares each model's births, deaths, and total population to a design with options that most closely replicate UN totals over a 150-year period.

## 1.2   Demographic Statistics from the United Nations

The World Population Prospect reports from the UN provide fertility, mortality, and migration statistics for all countries covering the year 1950–2100 based on a cohort component method (CCM) of population estimation and projection [3]. The UN's CCM from its 2019 report operates in 5-year steps and categorizes the population by 5-year age groups and by sex. It estimates and projects population totals for every 5-year interval using the following statistics: fertility rates, infant sex ratio, survival ratios, and net migration counts [4]. Understanding how to implement these statistics in a microsimulation requires specialized knowledge. The UN's statistics imply a specific event ordering and require more adjustments when adapting to a 1-year step simulation. Our team developed a microsimulation design called "Split Fertility", which closely replicates the UN totals in both a 5- and 1-year time step, by carefully interpreting the mathematical relationships and assumptions latent in the UN's CCM approach.

## 1.3   The "Split Fertility" Design

The Split Fertility (SF) design adapts the UN's CCM approach to a microsimulation [5]. It specifies a particular ordering of demographic events where fertility occurs twice during every time step, as follows: fertility round one, mortality, aging, migration, fertility round two. This order preserves the mathematical relationships of each statistic while operating prospectively within a simulation. It replicates UN data very closely when operating in 5-year steps.

In the 1-year version of the SF design, some statistics assign risk according to birth cohort instead of age. The fertility rate can apply to female agents' current age in each year, but the survival ratio and the migration count follow the birth cohort, which is labeled according to their ages at the beginning or end of the 5-year interval, respectively. In addition to tracking agents' current age, we also track their birth cohort and use this to assign mortality and migration events. Another small adjustment is to make immigrants skip the mortality event until the next 5-year interval starts, at which point they behave like non-migrants. This is done because

the mortality event applies to the cohort present at the start of the 5-year interval, which explicitly excludes immigrants.

These model options are selected because they replicate the implied event ordering and assignment of risk dictated by the UN's CCM and demographic statistics. The SF design produces under 1% divergence in total population, births, and deaths over a 150-year period when compared to the UN targets in the United States, India, and Norway [5]. This is true when we initialize the model with 100,000 agents and use the sorting method [1] to reduce stochasticity in demographic events. Details of the SF design are in Bacon et al. [5].

## 1.4  Alternative Option Combinations

The SF design is accurate to UN targets over a long span of time and in countries of different sizes and population dynamics, but its specific option combinations are not intuitive or well-known, and may not be appropriate for all project needs. Projects that prioritize efficiency, realism, or a more conventional microsimulation design in general may select different combinations of options. For example, initializing with a sample of 50,000 agents is very inefficient in some projects, and the convention of stochasticity may be preferred over variance-reduction techniques. It is also likely that someone can download the UN data without the knowledge of best practices for implementing UN demographic statistics in a microsimulation.

We compare the SF design to different combinations of options, both general options and those more specific to the data source. This highlights the importance of the specific SF design to adapt external demographic statistics to a microsimulation, but also identifies seemingly mundane options as very important for maintaining accurate population dynamics. General options that may be selected based on efficiency, convention, or project needs include the time step of the model (5 vs. 1 year steps), initial sample sizes (ranging from 500 to 50,000 agents), and whether to permit stochasticity in setting the age/sex distribution at initialization and in agents' risk of experiencing fertility and mortality events (migration events should not be treated as stochastic because the statistic is based on a count instead of a probability).

Options more specific to the data source include the order of demographic events. We contrast the "Split Fertility" ordering, described in Sect. 1.3, with a simpler fixed order where fertility occurs only once at the beginning, followed by mortality, migration, and finally aging. Those unfamiliar with the UN statistics might find this order reasonable and intuitive, since it first allows new agents to be born, all agents then experience mortality, migrants arrive/leave, then everyone ages before the next time step. Since the UN's mortality statistic applies to those present at the start of each 5-year interval, it is sensible for immigrants to skip death events until the next 5-year interval. Those unfamiliar with the statistic may not know this, so we include immigrant "immunity" as one design option. Lastly, the UN's mortality and migration statistics apply to birth cohort designations, which is a fixed characteristic,

rather than current age which updates with the aging event. Using cohort instead of current age for these events is another data-specific design option.

## 2   Methods

### 2.1   Data

Data for simulation inputs and validation come from the UN's 2019 World Population Prospects report for three countries: Norway, the United States (US), and India [3]. These three countries vary tremendously in population growth over the 150 years, and differ in population dynamics as well, which challenges the model specifications. Norway and the USA have experienced exceptional immigration, particularly in recent decades, but the USA has grown more quickly than Norway. India differs from both, in that it has very low migration, but much higher fertility, mortality, and overall population growth, particularly in the twentieth century.

### 2.2   Model Design Options

There are 176 unique model designs; 48 operate in 5-year steps and 128 in 1-year steps. Each model is run separately by country, resulting in a total of 528 models. We designed all microsimulations using AnyLogic 8.7.6., where each model is replicated 30 times to achieve high confidence in results that vary due to stochasticity. There are seven option toggles, two of which are conditional on the absence of another option. These are as follows:

1. **The time step of the model**: 5-year steps versus 1-year steps. Common-sense adjustments adapt the UN 5-year statistics to annual steps. Individual project needs determine which time step is preferred.
2. **Agent sample size at initialization**: 500, 1000, 10,000, 50,000. Many demographic microsimulations use a set percentage of the population (e.g. 1 or 2%) [6]. In very large countries (e.g. India), even a small starting percentage is a large number of agents and may be inefficient. We scale down each country's population to the same number for comparability.
3. **Variance reduction at initialization**: Whether stochasticity is reduced in setting agents' age/sex distribution at initialization. With stochasticity, agents' age/sex characteristics are set using a probability based on the UN's age/sex proportions for each country in 1950. If variance reduction is used, agents' age/sex characteristics exactly match the UN's population distributions in 1950 each time the model is run.
4. **Variance reduction in fertility/mortality events**: Whether stochasticity is reduced in agents' risk of fertility and mortality events. With stochasticity, agents'

risk of fertility/mortality are set using the UN's age/sex-specific statistics as probabilities. If the sorting method for variance reduction is used, the UN statistic is applied to the population of agents at risk to identify the same number of affected agents each time the model is run.

5. **Event ordering**: Split Fertility fixed event ordering vs. a simpler fixed order. The Split Fertility order has fertility occur twice with a halved probability in each occurrence and before/after all other events (fertility, mortality, aging, migration, fertility). This contrasts to a simpler order with fertility, mortality, migration, and aging as the set order.

6. **Immigrant immunity to death** (excluding 5-year models): Whether immigrant agents skip the mortality event. This option is redundant when operating in 5-year steps because all immigrants arrive after the mortality event in both event ordering options.

7. **Risk assignment by age and cohort** (excluding 5-year Split Fertility ordering models): Whether current age or birth cohorts are used to assign mortality and migration events. This toggle is redundant when operating in 5-year steps with the Split Fertility event ordering, because current age and birth cohort designations always have the same value.

## 2.3 Analysis of Model Results

Analyses presented in tables use the *average* value calculated across a model design's 30 replications. Divergence is calculated as the difference between each model's results when compared to the UN's total expected totals in births, deaths, and population count across all time intervals. For the sake of brevity, we only present results on total population, suppressing information on divergences in births and deaths, which can be useful in determining the reasons for divergence in the total population. A value of 5% divergence in total population, for example, indicates that the number of agents exceeds or falls short of the UN population targets by 5% on average, when using a given model option. For ease of interpretation, we identify what percent of models with the design option are "top performing" models (have less than 5% divergence in total population). We also present line graphs of individual model results, which highlights the variation in magnitude and directionality of the divergences over time and in each country context.

## 3 Results

Table 1 shows the overall prevalence of top-performing models in Norway, the US, and India, and then by model design option. Results are separated by 5-year and 1-year steps. Among the 5-year models, about 73% of the Norway models have low divergence while only 13% of India models do, but the difference by country in the

1-year models is much smaller. This indicates that many model designs result in high divergence in countries with high growth rates, and the 1-year models are particularly sensitive to model design options.

In the 5-year results, models with at least 1,000 agents can achieve under 5% divergence, but none of the models with only 500 agents have low divergence. In the 1-year models, the threshold for sufficient agent sample is much higher; most of the designs with fewer than 50,000 agents have higher divergence. About half of models achieved low divergence when using variance reduction at initialization and fertility/mortality events. This was also the case for the event ordering option. Although only half of 5-year models with Split Fertility ordering and the cohort risk assignment are top-performing, about 70% of models using cohort risk in the 1-year models have low divergence. Whether immigrants skip the death event appears to make little difference. This is likely because their immunity only lasts for the first five years of residence and countries with large immigrant flows have low mortality rates.

Next, the results of individual models are displayed in graphs for each country (Fig. 1). Graphed results show the *best-performing* replication of each model design. Total overall population is presented separately for each country, and black dots indicate the UN "target" values at each 5-year interval. The 5-year designs are in blue and the 1-year designs in red. The shading gradient indicates variation in sample size, where the darkest indicates a larger initial agent population (e.g. 50,000) and

**Table 1**  Percent of models achieving below 5% divergence in total population

|  | 5-year models | | 1-year models | |
|---|---|---|---|---|
|  | % | N | % | N |
| All models from Norway, the United States, and India | 38 | 144 | 26 | 384 |
| Norway | 73 | 48 | 33 | 128 |
| United States | 29 | 48 | 28 | 128 |
| India | 13 | 48 | 17 | 128 |
| *Number of agents in initial sample* | | | | |
| 50,000 | 33 | 36 | 76 | 96 |
| 10,000 | 33 | 36 | 22 | 96 |
| 1,000 | 33 | 36 | 2 | 96 |
| 500 | 0 | 36 | 0 | 96 |
| Stochastic age/sex initialization | 45 | 72 | 48 | 192 |
| Stochastic fertility and mortality events | 42 | 72 | 53 | 192 |
| Split fertility event order | 56 | 72 | 44 | 192 |
| Cohort age[a] | 50 | 48 | 68 | 192 |
| Immigrants skip mortality event[b] | NA | 0 | 50 | 192 |

[a] Option only present in 5-year models and in 1-year split fertility event order models
[b] Option only present in 5-year models

lightest is the smallest (e.g. 500). All displayed values have been scaled up to true-population sizes, for comparability to the expected dynamics in each country. As shown, 1-year steps noticeably produce results with more extreme divergence, and this is particularly the case when initial samples are smaller. Some 1-year designs even result in extinctions.

Table 1 and Fig. 1 indicate that time step and sample size are of primary importance. What value, then, does the specific "SF" design [5] add? The SF design achieves very low divergence in countries with different population dynamics, but this accuracy is contingent on sample size. Figure 2 shows results by time step and initial sample size, with the Split Fertility specification as follows: "SF" event order, variance reduction at initialization and for mortality/fertility events, cohort risk assignment, and immigrant death immunity. As shown, the 5-year models match the UN target over the entire time span with at least 1,000 agents, and similar accuracy is achieved in the 1-year model with 50,000 agents.

It is possible to achieve similarly low divergence in a small sample of agents with other design combinations, but at a cost. The specific design combinations are unpredictable, and accuracy of the model may be unsuitable for use in other projects.



**Fig. 1** Total population from each model design's best performing replication in 5-year steps (blue) and 1-year steps (red) against the UN targets (black dots)

**Fig. 2** Total population using the Split Fertility model design in 5-year steps (blue) and 1-year steps (red) against the UN targets (black dots)

For example, consider a design with only 500 or 1000 agents that operates in 1-year steps, has variance reduction for agents' age/sex initialization but not fertility/mortality events, uses a simple event order (not Split Fertility), and uses cohort age to assign risk. Results from this design are shown in Fig. 3. This unique combination of options interacts with India's population dynamics to achieve total population divergence of 6% or less. Using the same design in Norway or the US, however, leads to substantial divergence (Fig. 3.). Developing a microsimulation without regard to the mathematical relationships between the UN statistics could produce a misleading level of fidelity in one country, and perform poorly if implemented in another country context.

## 4   Discussion and Conclusions

When designing a microsimulation, options may be selected for efficiency, realism, accuracy, or simply to follow convention in lieu of certainty on the appropriate design. Our model option comparisons illustrate how certain combinations of options replicate expected UN population dynamics reliably in different countries, and

**Fig. 3** Total population using a small sample model design in 5-year steps (blue) and 1-year steps (red) against the UN targets (black dots)

that varying time step and the initial sample size of agents are among the most consequential options.

Our results show that small samples produce high divergence, particularly in 1-year step models. In small samples and shorter time steps, event probabilities become increasingly rare, even for seemingly common events such as giving birth. This is demonstrated in the US and Norway; as fertility rates decline some models produce no births. This problem is less severe in India, because initially high fertility and rapid population growth increase the number of agents in the model very quickly, making it possible for rare events to affect at least one agent in subsequent time intervals. The need for large samples when using variance reduction techniques like such as sorting method is already known, but permitting stochasticity has its own cost and cannot solve the problem in very small samples either [1].

The results also generally support the Split Fertility (SF) design as reducing divergence, but it cannot guarantee the success of SF designs in smaller initialized samples, and could even make divergence worse in very small sample populations because halving the fertility event's probability makes it a rarer event. Although the specific SF design is ideal in handling the mathematical relationships in the UN's demographic statistics in multiple countries, it requires a large initial sample of agents when operating in 1-year steps, which may be prohibitive for some projects. We also identified

the possibility of small sample designs that can replicate UN targets in India, but at the cost of predictability and generalizability to other country contexts. Knowledge of these drawbacks can aid in the design and validation of microsimulations that use UN data or similar demographic statistics to govern population dynamics.

Our research demonstrates the importance of certain assumptions and model designs, and also highlights how seemingly mundane design options can make the adaptation of external statistics surprisingly difficult. Our team's initial troubleshooting of the SF model design, for instance, required many hours of testing models with large samples of agents. These kinds of tests are not always feasible. We provide the AnyLogic and R code for this paper on a GitHub repository accessible at: https://github.com/centerformindandculture/UN-CCM-1-MICROSIM. We invite others to re-use our code under a Create Commons license to assess the effect of selecting different model options in their own demographic microsimulations. More research is needed to provide guidance on the development and validation of microsimulations.

# References

1. Van Imhoff, E., Post, W.: Microsimulation methods for population projection. In: Population: An English Selection, 97–138 (1998)
2. Müller, K., Axhausen, K.W.: Population synthesis for microsimulation: state of the art. Arbeitsberichte Verkehrs- Und Raumplanung, 638 (2010)
3. United Nations. World Population Prospects 2019: Methodology of the United Nations Population Estimates and Projections (2019)
4. Marois, G., KC, S.: Converting a cohort component model into a microsimulation model. In: Marois, G., KC, S. (eds.) Microsimulation Population Projections with SAS: A Reference Guide, pp. 25–49 (2021)
5. Bacon, R., Hodulik, G., Voas, D., Puga-Gonzalez, I., Wildman, W.: Stay on target: population projections and microsimulation design. Int. J. Microsimulation. **16**(3), 1–19 (2023)
6. Li, J., O'Donoghue, C.: A survey of dynamic microsimulation models: uses, model structure and methodology. Int. J. Microsimulation. **6**(2), 3–55 (2013)

# On Social Simulation in 4D Relativistic Spacetime

**Lai Kwun Hang**

**Abstract**  Agent-based modeling and simulation allow us to study social phenomena in hypothetical scenarios. If we stretch our imagination, one of the interesting scenarios would be our interstellar future. To model an interstellar society, we need to consider relativistic physics, which is not straightforward to implement in existing agent-based simulation frameworks. In this paper, we present the mathematics and algorithmic details needed for simulating agent-based models in 4D relativistic spacetime. These algorithms form the basis of our open-source computational framework, "Relativitization" (Lai in [1]).

**Keywords**  Special relativity · Interstellar society · Agent-based simulation · Computational framework

## 1   Introduction

The scientific and technological advancement in the last century greatly increases our understanding of the universe. Nowadays, we are able to build giant telescopes and observe astronomical objects billions of light years away. Apart from deepening our scientific understanding, our astronomical knowledge also stimulates our imagination of interstellar civilizations. A lot of great science fiction has been written, and scientists have proposed ideas like the Fermi paradox [2], the Dyson sphere [3] and the Kardashev scale [4]. While many of these ideas are physically plausible, it would be interesting to discuss these ideas from a social science perspective. Due to the highly hypothetical nature of the problem, we suggest that agent-based model (ABM) enables formal academic discussion on interstellar society.

An ABM is a simulation model bridging microscopic behaviours of agents and macroscopic observations. Depending on the context of the model, an agent can be an individual, an organization, or even a country. First, assumption about the social behaviours of the agents are made, then the agents are placed and evolved in a com-

L. K. Hang (✉)

Centre for Science and Technology Studies (CWTS), Leiden University, Leiden, Netherlands
e-mail: k.h.lai@cwts.leidenuniv.nl

putational environment. For a research problem for which it is not viable to perform data collection and analysis, ABM can still be used for theoretical exposition [5].

To model agents in interstellar space, suppose we only consider a scale with normal stellar objects so that we can ignore the effects of general relativity, such as universe expansion and black holes, we still have to consider the effect of special relativity. In the context of ABM, where the simulation is computed under a set of inertial frames that are at rest to each other, we can simplify the theory of relativity into two core phenomena: speed of light as the upper bound of the speed of information travel, and time dilation relative to any stationary observer in the inertial frames. This also implies that we have to take care of four dimensions: three space dimensions, plus one time dimension.

Typically, an ABM is constructed using an ABM framework to facilitate model development and communication. There are a lot of existing ABM frameworks, to name a few, NetLogo [6], mesa [7], and Agents.jl [8], see [9] for a detailed review. While it is possible to build a 4D relativistic model in some existing ABM frameworks, those frameworks do not have native support for the necessary 4D data structures, and it can be error-prone to enforce relativistic effects via custom implementations of data structures and algorithms. Therefore, we have developed a simulation framework we call "Relativitization" [1], to help social scientists to build an ABM in relativistic spacetime. In this paper, the mathematics and the algorithms underlying the framework will be presented.

## 2   Definitions

In Relativitization, an agent is called a "player". Players live in a universe. Ideally, computation should be done in every local frame following all players, and the computation results can be synchronized by Lorentz transformations. However, this will make the framework and the model substantially more complex. Therefore, all computations are done according to some inertial frames that are at rest to each other. The spatial coordinates of a player are represented by floating-point numbers $x$, $y$ and $z$ and the time coordinate of a player is represented by a floating-point number $t$. To simplify computation and visualization, the universe is partitioned into unit cubes. A player with floating-point coordinates $(t, x, y, z)$ is located at the cube with integer coordinates $T = \lfloor t \rfloor$, $X = \lfloor x \rfloor$, $Y = \lfloor y \rfloor$, $Z = \lfloor z \rfloor$, note that the computations of a simulation are done at unit time steps and we can actually assume $T = t$. Denote the speed of light as $c$. In vector notation, define $\mathbf{s} = (t, \overrightarrow{u}) = (t, x, y, z)$, and $\mathbf{S} = (T, \overrightarrow{U}) = (T, X, Y, Z)$.

## 2.1 Interval and Time Delay

The spacetime interval between coordinates $\mathbf{s}_i$ and $\mathbf{s}_j$ is

$$\|\mathbf{s}_i - \mathbf{s}_j\| = c^2(t_i - t_j)^2 - (x_i - x_j)^2 - (y_i - y_j)^2 - (z_i - z_j)^2. \tag{1}$$

If $\|\mathbf{s}_i - \mathbf{s}_j\| < 0$, it is called a spacelike interval, and events that happen at the two coordinates are not causally connected because no information can travel faster than the speed of light $c$.

It is often needed to compute intervals in integer coordinates. We define the spatial distance between $\overrightarrow{U_i}$ and $\overrightarrow{U_j}$ as the maximum distance between all points in the cubes at $\overrightarrow{U_i}$ and $\overrightarrow{U_j}$

$$|\overrightarrow{U_i} - \overrightarrow{U_j}| = (X_i - X_j + 1)^2 + (Y_i - Y_j + 1)^2 + (Z_i - Z_j + 1)^2. \tag{2}$$

Suppose there is a signal sent from $\overrightarrow{U_i}$ to $\overrightarrow{U_j}$. To ensure that the information travels slower than the speed of light, the integer time delay $\tau(\overrightarrow{U_i}, \overrightarrow{U_j})$ is computed as

$$\tau(\overrightarrow{U_i}, \overrightarrow{U_j}) = \left\lceil \frac{|\overrightarrow{U_i} - \overrightarrow{U_j}|}{c} \right\rceil. \tag{3}$$

## 2.2 Group Id

From Eq. 2, even if $\overrightarrow{U_i} = \overrightarrow{U_j}$, the time delay is non-zero. To implement zero time delay for players that are really close to each other, we divide a unit cube into several sub-cubes with edge length $d_e$, and information travel within the same sub-cubes is instantaneous.

To improve the computational speed when checking whether two players belong to the same sub-cube, we assign a "group id" to each sub-cube in a unit cube. A unit cube has $n_e^3$ sub-cubes, where $n_e = \left\lceil \frac{1}{d_e} \right\rceil$. For a player at $\overrightarrow{u}$, it belongs to the $(n_x, n_y, n_z)$ sub-cubes, where $n_x = \left\lfloor \frac{x-X}{d_e} \right\rfloor$, $n_y = \left\lfloor \frac{y-Y}{d_e} \right\rfloor$, and $n_z = \left\lfloor \frac{z-Z}{d_e} \right\rfloor$. The group id $g(\overrightarrow{u}, \overrightarrow{U})$ of the player can be computed as

$$g(\overrightarrow{u}, \overrightarrow{U}) = n_x n_e^2 + n_y n_e + n_z. \tag{4}$$

If two players have the same integer coordinates and the same group id, then we say the players belong the same group and the time delays between the players are zero.

## 2.3  Player Data

A player is characterized by a set of data:

- player id $i$,
- integer coordinates $(T_i, X_i, Y_i, Z_i)$,
- a historical record of integer coordinates $H_i = \{(T_i', X_i', Y_i', Z_i') \mid T_i' < T_i\}$,
- floating-point coordinates $(t_i, x_i, y_i, z_i)$,
- time dilation counter variables $\mu_i$, a floating point variable, and $\nu_i$, a boolean variable, to keep track of time dilation (see Sects. 2.6 and 3.3), $\mu_i = 0$ initially for all players,
- group id $g_i$,
- floating-point velocities $\overrightarrow{v_i} = (v_{ix}, v_{iy}, v_{iz})$,
- other data $D_i$ relevant to the model.

## 2.4  Command

In other frameworks, interactions in ABMs are often presented as one player asking another player to do something. Because the speed of information travel is bounded by $c$, a player cannot simply ask other players to do something immediately. Instead, interactions are mediated by commands. Whenever player $i$ wants to interact with player $j$, player $i$ sends a command to $j$.

A command is characterized by:

- $i_{\text{to}}$, the id of the player to receive this command,
- $i_{\text{from}}$ the id of the target player who sent this command,
- $\mathbf{S}_{\text{from}}$ the integer coordinates when the player sent this command,
- $f_{\text{target}}$ a function to modify data of the target player when this is received.

Commands travel at the speed of light $c$. The amount of time needed for a command to reach the target, measured in the inertial frames we used in the simulation, depends on the trajectory of the target player $i_{\text{to}}$ and the sender coordinates $\mathbf{S}_{\text{from}}$.

## 2.5  Universe Data

Universe is an overarching structure which aggregates all necessary data and functionalities. An universe has:

- a current universe time $T_{\text{current}}$,
- a 4-dimension array of maps from player id to lists of player data $M_{TXYZ}$, so that the data of a player residing at $(T, X, Y, Z)$ is stored in the associated list, the "afterimages" of players are also stored in the corresponding list (Sec 3.5),

- a map $M_{\text{command}}$ from player id to lists of commands, such that a command in the list will be executed when the player receive the command,
- other universe global data $D_G$ relevant to the model.

## 2.6 Mechanism

Given an instance of a universe, the dynamics of players are based on predefined rules and the state of the universe observed by the players. In our framework, we call the rules mechanisms. A mechanism takes the state of the universe observed by a player, modifies the state of a player, and generates a list of commands to send to other players.

To ease the model development to account for the time dilation effect, we further divide mechanisms into two categories: regular mechanisms and dilated mechanisms. A regular mechanism is executed once per turn, while a dilated mechanism is executed once per multiple turns, adjusted for the time dilation of the player measured in the inertial frames we used in the simulation.

## 3 Simulation Step

The following are needed to define a model:

- the data structure of other player data $D_i$,
- the data structure of other universe global data $D_G$,
- a set of available commands,
- a function to initialize the universe data,
- a function to update the universe global data,
- a set of regular and dilated mechanisms.

Along with the universe data, it is useful to define a map $M_{\text{current}}$ from a player id to the current player data, i.e., $T_i = T_{\text{current}}$, as an internal object of the simulation. The modifications of player data are first performed on $M_{\text{current}}$, and then synchronized back to the universe data at appropriate timing.

Suppose we have initialized an universe model and $M_{\text{current}}$, a complete step in a simulation involves:

1. update the global data (Sect. 3.1),
2. compute time dilation effects for all players (Sect. 3.2),
3. process mechanisms for each player (Sect. 3.3),
4. process the command map (Sect. 3.4),
5. move players, add afterimages, and update time (Sect. 3.5).

The simulation can be ran for a fixed amount of steps, or stop when a stopping condition is met.

## 3.1 Update Global Data

A model may rely on a mutable global data $D_G$ to implement the dynamics. If the model depends on some player data to update $D_G$, and the effect is observable by players, we need to ensure that no information is transferred faster than the speed of light via the global data update.

For example, if the global data is modified if "all" player data satisfy a condition, we have to be careful about what we mean by "all" here. In the universe, the maximum time delay equals $\tau_{max} = \tau((0, 0, 0), (max(X), max(Y), max(Z)))$. To fulfill the speed of light constraint, the update function has to check whether all player data in $M_{TXYZ}$, where $T_{\text{current}} - \tau_{max} \leq T \leq T_{\text{current}}, 0 \leq X \leq max(X), 0 \leq Y \leq max(Y)$, and $0 \leq Z \leq max(Z)$, satisfy that condition.

## 3.2 Compute Time Dilation

Relative to a stationary observer $j$ in an inertial frame, special relativity predicts that a moving observer $i$ experiences a time dilation effect:

$$\gamma_i = \frac{1}{\sqrt{1 - \frac{v_i^2}{c^2}}}, \tag{5}$$

$$\Delta t_i = \frac{\Delta t_j}{\gamma_i}, \tag{6}$$

where $\gamma_i$ is called the Lorentz factor.

To account for the time-dilation effect, the time dilation counter variables $\mu_i$ and $\nu_i$ are updated by Algorithm 1 every turn for every player. $\nu_i$ will then affect the mechanism processing in Sect. 3.3.

---

**Algorithm 1:** Update time dilation counter

---

**Input**: $M_{\text{current}}$, map from player id to current player data

1 **foreach** player $i$ in $M_{\text{current}}$ **do**

2      $\mu_i \leftarrow \mu_i + \sqrt{1 - \frac{v_i^2}{c^2}}$;

3      **if** $\mu_i \geq 1$ **then**

4          $\mu_i \leftarrow \mu_i - 1$;

5          $\nu_i \leftarrow$ **true**;

6      **end**

7      **else**

8          $\nu_i \leftarrow$ **false**;

9      **end**

10 **end**

---

## 3.3   Process Mechanisms

Before processing any mechanism for a player, we need to compute the state of the universe viewed by the player. At an instance in our discretized relativistic universe, player $i$ sees other players located at the unit cubes closest to the surface of the past light cone of player $i$, while the entire cubes are still within the past light cone. The computation consists of two steps: (1) Algorithm 2 computes the view centered at a specific cube, ignoring the zero time delay when players are within the same group, (2) Algorithm 3 computes the view for players in a group. Assuming each line of these algorithms takes $O(1)$ and iterating over all $(X, Y, Z)$, the time complexity $O(mn)$ from Algorithm 2 dominates, where $m = X_{max} Y_{max} Z_{max}$ is the spatial size of the universe, and $n$ is the number of player.

---

**Algorithm 2:** Compute the view of the universe at a cube, ignore the zero time delay when player are in the same group

---

  **Input**:
    $T_i, X_i, Y_i, Z_i$ position of the viewing location
    $M_{TXYZ}$ 4D array of maps from player id to lists of player data
  **Output**:
    $M$ map from player id to player data
    $\Lambda_{XYZ}$ 3D array of maps from group id to lists of player id
1 Initialize empty $M$ and $\Lambda_{XYZ}$;
2 **foreach** $X_j, Y_j, Z_j$ **do**
3     $T_j \leftarrow T_i - \tau(\overrightarrow{U_i}, (X_j, Y_j, Z_j))$;
4     **foreach** player data in $M_{T_j X_j Y_j Z_j}[k]$ **do**
5        **if** $M$ has key $k$ **then**
6           **if** $T$ of $M[k] < T$ of the new player data **then**
7              Replace $M[k]$ by this new player data;
8           **end**
9        **end**
10        **else**
11           Store data of player $k$ to $M[k]$;
12        **end**
13     **end**
14 **end**
15 Associate the player id from $M$ to the corresponding list in $\Lambda_{XYZ}$ by spatial coordinates and group id;
16 **return** $(M, \Lambda_{XYZ})$

---

The view of the universe of a player is used by mechanisms to update the player data and generate commands to send. Regular mechanisms update the data of the player each turn, while dilated mechanisms update the player if the time dilation counter variable $v_i$ is **true** to account for the time dilation effect. The generated commands are executed immediately if the target player is within the same group of the sender, otherwise the commands are stored in $M_{command}$. Algorithm 4 shows the overall iterative process.

---

**Algorithm 3:** Compute the view of the universe for players in a group

---

**Input**:
    $g_i$ group id
    $T_j, X_j, Y_j, Z_j$ position of the viewing location
    $M$ map from player id to player data
    $\Lambda_{XYZ}$ 3D array of maps from group id to lists of player id
    $M_{TXYZ}$ 4D array of maps from player id to lists of player data
**Output**:
    $M'$ map from player id to player data
    $\Lambda'_{XYZ}$ 3D array of maps from group id to lists of player id
**1** $M' \leftarrow M$;
**2** $\Lambda'_{XYZ} \leftarrow \Lambda_{XYZ}$;
**3 foreach** player data in $M_{T_j X_j Y_j Z_j}[k]$ where $g(\vec{u_k}, \vec{U_k}) = g_i$ **do**
**4**    **if** $T$ of $M[k] < T$ of the new player data **then**
**5**       Replace $M'[k]$ by this new player data;
**6**       Update the corresponding position of player $k$ in $\Lambda'_{XYZ}$;
**7**    **end**
**8 end**
**9 return** $(M', \Lambda'_{XYZ})$

---

**Algorithm 4:** Update all players by mechanisms

---

**Input**:
    $M_{\text{current}}$ map from player id to current player data
    Universe data
**1 foreach** $(X_j, Y_j, Z_j)$ **do**
**2**    Compute the view of the universe at this cube by algorithm 2;
**3**    **foreach** group in this cube **do**
**4**       Compute the view of the universe at this group by algorithm 3;
**5**       **foreach** data of player $k$ in this group **do**
**6**          Update $M_{\text{current}}[k]$ by all regular mechanisms;
**7**          **if** $v_k$ is **true then**
**8**             Update $M_{\text{current}}[k]$ by all dilated mechanisms;
**9**          **end**
**10**       **end**
**11**       **foreach** generated command where target player $l$ is in this group **do**
**12**          Update $M_{\text{current}}[l]$ by $f_{\text{target}}$ of the command;
**13**       **end**
**14**    **end**
**15 end**
**16** Add the rest of commands to $M_{\text{command}}$ by the target player id of the commands;

### 3.4 Process Command Map

The command map $M_{command}$ is a map from player id to lists of commands that is being sent to that player. At each turn, the distance between the target player and the sent positions of all commands in the list are calculated, and the command is executed on the player if the spacetime interval is larger than zero. Algorithm 5 illustrates the process.

---

**Algorithm 5:** Process command map.

---

**Input**:
    $M_{current}$ map from player id to current player data
    $M_{command}$ map from player id to lists of commands

1 **foreach** key $i$ in $M_{command}$ **do**
2     Get the integer coordinates $\mathbf{S}_i$ of player $i$ from $M_{current}[i]$;
3     **foreach** command $C$ in the list $M_{command}[i]$ **do**
4         **if** $\|\mathbf{S}_{from} - \mathbf{S}_i\| \geq 0$ **then**
5             Update $M_{current}[i]$ by $f_{target}$ of the command;
6         **end**
7     **end**
8 **end**
9 Remove all executed commands;

---

### 3.5 Move Players and Add Afterimages

Moving players and storing their data requires additional considerations in this simulation framework. Consider the following example:

1. assume player $i$ and player $j$ are located in the same cube,
2. player $j$ moves to the other cube,
3. the new information takes time to travel to player $i$, so player $i$ cannot see the new position of player $j$,
4. player $i$ cannot see the old information of player $j$ either, because player $j$ is no longer there,
5. player $j$ disappears from the sight of player $i$.

This "disappearance" is caused by the problem of the integer-based coordinates used in the computation of player's 3D view.

Consider a more generic situation: suppose player $i$ is located at $\overrightarrow{U_i}$, and player $j$ moves from $\overrightarrow{U_j}$ to $\overrightarrow{U_k}$. Ignoring the possibility of zero time delay, the maximum time player $i$ has to wait to see player $j$ is bounded by Eq. 7,

$$\Delta T = \tau(\overrightarrow{U_i}, \overrightarrow{U_j}) - \tau(\overrightarrow{U_i}, \overrightarrow{U_k}), \tag{7}$$

$$= \left\lceil \frac{|\overrightarrow{U_i} - \overrightarrow{U_j}|}{c} \right\rceil - \left\lceil \frac{|\overrightarrow{U_i} - \overrightarrow{U_k}|}{c} \right\rceil, \tag{8}$$

$$\leq \left\lceil \frac{|\overrightarrow{U_i} - \overrightarrow{U_j}|}{c} - \frac{|\overrightarrow{U_i} - \overrightarrow{U_k}|}{c} \right\rceil, \tag{9}$$

$$\leq \left\lceil \frac{|\overrightarrow{U_j} - \overrightarrow{U_k}|}{c} \right\rceil, \tag{10}$$

$$= \tau(\overrightarrow{U_j}, \overrightarrow{U_k}). \tag{11}$$

Therefore, if we include back the possibility where the time delay between player $i$ and player $j$ can be zero, the maximum duration of the disappearance produced by the movement is bounded by $\Delta T_{max} = \tau((0, 0, 0), (1, 1, 1))$. To prevent the unrealistic disappearance from happening, the old player data has to stay at the original position for at least $\Delta T_{max}$ turn, we call this the "afterimage" of the player. Note that afterimages only participate in the 3D view of players, they should not be updated by commands or mechanisms.

Algorithm 6 does multiple things: it updates the universe time, it moves players by their velocities, it synchronizes time of players, it stores old coordinates to the history of player, it cleans the history if the stored coordinates is too old, and it adds the current player and afterimages to the latest spatial 3D array in the 4D data array $M_{TXYZ}$. Since the universe time has been updated, this simulation step has finished, the universe should go to the next step and loop over all algorithms in Sect. 3 again.

## 4   Discussion

The presented algorithms form the backbone of our computational framework, "Relativitization" [1]. There are technical subtleties that are not discussed here, such as creating new players, removing dead players, introducing randomness to models, parallelization of the algorithms, generating deterministic outcomes from parallelized simulations with random number generators, interactive human input to intervene in a simulation, etc. Nevertheless, the framework implements the major part of the technical subtleties, and provides a suitable interface to ease the development of any 4D, relativistic ABM.

It can be interesting to implement a classical ABM into the framework. Spatial ABMs with non-local interactions, such as the classical flocking model [10], are particularly suitable. These models are naturally affected by the time delay imposed by the speed of light limitation. Simulating such a model in the Relativitization framework allows us to explore the effects of time delay on the model.

---

**Algorithm 6:** Move player and add afterimages

---

**Input**:

$M_{current}$ map from player id to current player data

Universe data

1 $T_{current} \leftarrow T_{current} + 1$;

2 Initialize a 3D array of maps from player id to lists of player data $M_{XYZ}$;

3 **foreach** data of player $i$ in $M_{current}$ **do**

4      $t_i \leftarrow T_{current}$;

5      $x_i \leftarrow x_i + v_{ix}$;

6      $y_i \leftarrow y_i + v_{iy}$;

7      $z_i \leftarrow z_i + v_{iz}$;

8      $T_i \leftarrow T_{current}$;

9      $X_i \leftarrow \lfloor x_i \rfloor$;

10      $Y_i \leftarrow \lfloor y_i \rfloor$;

11      $Z_i \leftarrow \lfloor z_i \rfloor$;

12      $g_i \leftarrow g(\vec{u_i}, \vec{U_i})$ by Eq. 4;

13      **if** coordinates or group is new **then**

14          Save the previous coordinates to history $H_i$;

15      **end**

16      **foreach** $(T'_i, X'_i, Y'_i, Z'_i)$ in $H_i$ **do**

17          Remove from $H_i$ if $T_{current} - T' > \Delta T_{max}$;

18      **end**

19      Save the new data to $M_{X_i Y_i Z_i}[i]$;

20      **foreach** $(T'_i, X'_i, Y'_i, Z'_i)$ in $H_i$ **do**

21          Find the old player data from $M_{T'_i X'_i Y'_i Z'_i}[i']$;

22          Add the old player data to $M_{X'_i Y'_i Z'_i}[i']$;

23      **end**

24 **end**

25 Drop the oldest 3D spatial array from $M_{TXYZ}$;

26 Add $M_{XYZ}$ as the latest spatial array to $M_{TXYZ}$;

---

Ultimately, existing ABMs might not be suitable to describe interstellar society. A solid understanding of social mechanisms and physics, together with some artistic imagination, are needed to build inspiring interstellar ABMs. As a first step, we have integrated a few social mechanisms to build a big "model", which is also a game. The "model" can be found on the GitHub[1] repository of our framework.

Apart from the possibility of implementing different models using the framework, the algorithms may also be optimized further. For example, the iteration in Sect. 3.3 has a time complexity of $O(mn)$. A naive alternative implementation to iterate over all the combinations of players could change the complexity to $O(n^2)$, which could have better performance when the density of players is low. We leave these potential improvements to future research.

---

[1] https://github.com/Adriankhl/relativitization.

## 5   Conclusion

In this paper, we have presented a set of algorithms to implement ABM simulations in a 4D, relativistic spacetime. Based on these algorithms, we have developed a simulation framework we call "Relativitization" [1]. Our framework will lower the barrier of entry for social scientists to apply their expertise to explore the interstellar future of human civilization. We hope our framework can be used to initiate meaningful and academically interesting discussions about our future.

## References

1. Lai, K.H.: Relativitization https://github.com/Adriankhl/relativitization, https://doi.org/10.5281/zenodo.6120765 (2022)
2. Gray, R.H.: The Fermi paradox is neither Fermi's nor a paradox. Astrobiology **15**, 195–199 (2015)
3. Wright, J.T.: Dyson spheres. arXiv preprint arXiv:2006.16734 (2020)
4. Gray, R.H.: The extended kardashev scale. Astronom. J. **159**, 228 (2020)
5. Edmonds, B., Meyer, R.: Simulating Social Complexity, Springer (2015)
6. Wilensky, U.: NetLogo http://ccl.northwestern.edu/netlogo/ (1999)
7. Kazil, J., Masad, D., Crooks, A.: Utilizing python for agent-based modeling: the mesa framework. In: Thomson, R., Bisgin, H., Dancy, C., Hyder, A., Hussain, M. (eds.). Social, Cultural, and Behavioral Modeling, pp. 308–317. Springer International Publishing, Cham (2020). ISBN: 978-3-030-61255-9
8. Datseris, G., Vahdati, A.R., DuBois, T.C.: Agents.jl: a performant and feature-full agent-based modeling software of minimal code complexity. SIMULATION **0**, 00375497211068820. eprint: https://doi.org/10.1177/00375497211068820, https://doi.org/10.1177/00375497211068820 (2021)
9. Pal, C.-V., Leon, F., Paprzycki, M., Ganzha, M.: A review of platforms for the development of agent systems. arXiv preprint arXiv:2007.08961 (2020)
10. Reynolds, C.W.: Flocks, herds and schools: a distributed behavioral model. In: Proceedings of the 14th Annual Conference on Computer graphics and Interactive Techniques, pp. 25–34 (1987)

# Author Index