

A detailed technical line drawing of a robotic arm, showing various joints, actuators, and sensors. The drawing is rendered in purple and white lines. The top portion of the image is overlaid with a solid orange background, which serves as a backdrop for the main title. The bottom portion of the image shows the full extent of the robotic arm's structure, including its base and various mechanical components.

COMPLEXITY AND SELF-ORGANIZATION

EDITED BY: Carlos Gershenson, Daniel Polani and Georg Martius
PUBLISHED IN: Frontiers in Robotics and AI, Frontiers in Physics and
Frontiers in Applied Mathematics and Statistics



frontiers

Frontiers eBook Copyright Statement

The copyright in the text of individual articles in this eBook is the property of their respective authors or their respective institutions or funders. The copyright in graphics and images within each article may be subject to copyright of other parties. In both cases this is subject to a license granted to Frontiers.

The compilation of articles constituting this eBook is the property of Frontiers.

Each article within this eBook, and the eBook itself, are published under the most recent version of the Creative Commons CC-BY licence.

The version current at the date of publication of this eBook is CC-BY 4.0. If the CC-BY licence is updated, the licence granted by Frontiers is automatically updated to the new version.

When exercising any right under the CC-BY licence, Frontiers must be attributed as the original publisher of the article or eBook, as applicable.

Authors have the responsibility of ensuring that any graphics or other materials which are the property of others may be included in the CC-BY licence, but this should be checked before relying on the CC-BY licence to reproduce those materials. Any copyright notices relating to those materials must be complied with.

Copyright and source acknowledgement notices may not be removed and must be displayed in any copy, derivative work or partial copy which includes the elements in question.

All copyright, and all rights therein, are protected by national and international copyright laws. The above represents a summary only. For further information please read Frontiers' Conditions for Website Use and Copyright Statement, and the applicable CC-BY licence.

ISSN 1664-8714
ISBN 978-2-88966-781-9
DOI 10.3389/978-2-88966-781-9

About Frontiers

Frontiers is more than just an open-access publisher of scholarly articles: it is a pioneering approach to the world of academia, radically improving the way scholarly research is managed. The grand vision of Frontiers is a world where all people have an equal opportunity to seek, share and generate knowledge. Frontiers provides immediate and permanent online open access to all its publications, but this alone is not enough to realize our grand goals.

Frontiers Journal Series

The Frontiers Journal Series is a multi-tier and interdisciplinary set of open-access, online journals, promising a paradigm shift from the current review, selection and dissemination processes in academic publishing. All Frontiers journals are driven by researchers for researchers; therefore, they constitute a service to the scholarly community. At the same time, the Frontiers Journal Series operates on a revolutionary invention, the tiered publishing system, initially addressing specific communities of scholars, and gradually climbing up to broader public understanding, thus serving the interests of the lay society, too.

Dedication to Quality

Each Frontiers article is a landmark of the highest quality, thanks to genuinely collaborative interactions between authors and review editors, who include some of the world's best academicians. Research must be certified by peers before entering a stream of knowledge that may eventually reach the public - and shape society; therefore, Frontiers only applies the most rigorous and unbiased reviews. Frontiers revolutionizes research publishing by freely delivering the most outstanding research, evaluated with no bias from both the academic and social point of view. By applying the most advanced information technologies, Frontiers is catapulting scholarly publishing into a new generation.

What are Frontiers Research Topics?

Frontiers Research Topics are very popular trademarks of the Frontiers Journals Series: they are collections of at least ten articles, all centered on a particular subject. With their unique mix of varied contributions from Original Research to Review Articles, Frontiers Research Topics unify the most influential researchers, the latest key findings and historical advances in a hot research area! Find out more on how to host your own Frontiers Research Topic or contribute to one as an author by contacting the Frontiers Editorial Office: frontiersin.org/about/contact

COMPLEXITY AND SELF-ORGANIZATION

Topic Editors:

Carlos Gershenson, National Autonomous University of Mexico, Mexico

Daniel Polani, University of Hertfordshire, United Kingdom

Georg Martius, Max Planck Institute for Intelligent Systems, Germany

Citation: Gershenson, C., Polani, D., Martius, G., eds. (2021). Complexity and Self-Organization. Lausanne: Frontiers Media SA.

doi: 10.3389/978-2-88966-781-9

Table of Contents

- 04** *Editorial: Complexity and Self-Organization*
Carlos Gershenson, Daniel Polani and Georg Martius
- 06** *Dynamical Inference of Simple Heteroclinic Networks*
Maximilian Voit and Hildegard Meyer-Ortmanns
- 16** *Metrics of Emergence, Self-Organization, and Complexity for EWOM Research*
Juan C. Correa
- 22** *How Computation is Helping Unravel the Dynamics of Morphogenesis*
David Pastor-Escuredo and Juan C. del Álamo
- 34** *Unsupervised Learning Facilitates Neural Coordination Across the Functional Clusters of the C. elegans Connectome*
Alejandro Morales and Tom Froese
- 41** *New Methods for the Steady-State Analysis of Complex Agent-Based Models*
Chico Q. Camargo
- 49** *How Could Future AI Help Tackle Global Complex Problems?*
Anne-Marie Grisogono
- 58** *Improving the Robustness of Online Social Networks: A Simulation Approach of Network Interventions*
Giona Casiraghi and Frank Schweitzer
- 69** *Collective Computation in Animal Fission-Fusion Dynamics*
Gabriel Ramos-Fernandez, Sandra E. Smith Aguilar, David C. Krakauer and Jessica C. Flack
- 84** *Computational Intelligence for Studying Sustainability Challenges: Tools and Methods for Dealing With Deep Uncertainty and Complexity*
Edmundo Molina-Perez, Oscar A. Esquivel-Flores and Hilda Zamora-Maldonado
- 102** *Personogenesis Through Imitating Human Behavior in a Humanoid Robot "Alter3"*
Atsushi Masumori, Norihiro Maruyama and Takashi Ikegami



Editorial: Complexity and Self-Organization

Carlos Gershenson^{1,2,3*}, Daniel Polani⁴ and Georg Martius⁵

¹ Departamento de Ciencias de la Computación, Instituto de Investigaciones en Matemáticas Aplicadas y en Sistemas, Universidad Nacional Autónoma de México, México, México, ² Centro de Ciencias de la Complejidad, Universidad Nacional Autónoma de México, México, México, ³ Lakeside Labs GmbH, Klagenfurt, Austria, ⁴ School of Engineering and Computer Science, University of Hertfordshire, Hatfield, United Kingdom, ⁵ Autonomous Learning Group, Max Planck Institute for Intelligent Systems, Tübingen, Germany

Keywords: complexity, self-organization, formal methods, cognitive science, global issues

Editorial on the Research Topic

Complexity and Self-Organization

Complexity occurs when relevant interactions prevent the study of elements of a system in isolation. These interactions between elements may lead to the self-organization of the system. A system can be described as self-organizing when its global properties are a product of the interactions of its components. Complexity and self-organization are prevalent in a broad variety of systems. Because of this, they have been studied from multiple perspectives and disciplines, leading naturally to transdisciplinary studies.

The scientific study of complexity and self-organization was limited before the popularization of computers in the 1980s, as previous tools were insufficient to deal with hundreds or thousands of variables. Thus, computer science has been essential for these studies.

In computational intelligence, complexity and self-organization have been studied and exploited with different purposes. The aim of this Research Topic was to bring together novel research into a coherent collection, spanning from theory and methods to simulations and applications.

For example, it has been observed that complex systems studied by different disciplines reach a balance between change and stability that has been also described as criticality. This balance allows the maintenance of functionality (robustness) and also the potential to change in response to new situations (adaptability). The different contributions included in this Research Topic illustrate how this balance is present in a broad variety of phenomena.

We received 22 submissions, out of which ten were accepted.

Of these, three dealt with **formal methods** as a crucial tool for the understanding of complex and self-organizing systems. Voit and Meyer-Ortmanns propose a data-based approach to automatically infer heteroclinic networks. These are networks where nodes represent saddle fixed points in phase space, while edges are heteroclinic orbits (which occur when the unstable manifold of a saddle fixed point intersects the stable manifold of another saddle). The proposed method is based on a template system that uses a learning algorithm to adjust the eigenvalues at the saddles, eventually reconstructing the topology of the original heteroclinic network. This approach is promising to infer a structure that reproduces an observed function or dynamics.

Camargo uses an agent-based model of ideological alignment to explore the usefulness of different approaches and methods to analyze its dynamics. Camargo generalizes to argue that the proposed approaches can be applied to other agent-based models of social behavior, where complexity is such that measurements of the performance of the model are not explicit nor straightforward.

An obstacle of measurements is the lack of proper metrics. In this respect, Correa presents a mini-review on how metrics of emergence, self-organization, and complexity can contribute to commerce/consumer studies, in particular, to the understanding “electronic word-of-mouth”

OPEN ACCESS

Edited and reviewed by:

Sheri Marina Markose,
University of Essex, United Kingdom

*Correspondence:

Carlos Gershenson
cgg@unam.mx

Specialty section:

This article was submitted to
Computational Intelligence in
Robotics,
a section of the journal
Frontiers in Robotics and AI

Received: 16 February 2021

Accepted: 02 March 2021

Published: 26 March 2021

Citation:

Gershenson C, Polani D and
Martius G (2021) Editorial: Complexity
and Self-Organization.
Front. Robot. AI 8:668305.
doi: 10.3389/frobt.2021.668305

(EWOM) data. These metrics can be used as proxies to the degree of customers' comments diversification, customers' comments polarization, and the diversification-polarization balance.

A particular subset of the studies of complex systems concerns itself with issues of **cognitive science**, which is addressed in three of the papers. The first of them asks what happens when humans interact with humanoid robots that mimic their behavior in a self-organizing way. This interesting question is investigated in Mazumori et al. In their system, the robot's behavior is partially self-organized using a memory of prior interaction with other humans (art gallery visitors) and an internal dynamics that makes the robot only partially predictable. Engagement of the interaction partner and the inversion of roles, i.e., that the human starts to imitate the robot, are frequently observed.

Ramos-Fernandez et al. consider the question how information is being acquired and distributed in a group of individuals, specifically validated against data obtained from spider monkey colonies. The particular dynamics studied are the dynamics of subgroups of monkeys which split and merge (fission-fusion dynamics) and how these decisions are taken by integrating events experienced by the individuals over time. The different timescales of the dynamics (fights, signaling, and rank) are used by the faster degrees of freedom as reference in a form of "downward causation." Thus, the decisions of the individuals of the collective make up the structure which, in turn, influences the decisions of the individuals.

Considering the neural mechanisms of behavior, Morales and Froese contribute to this research topic with a short study on the role of unsupervised learning in the formation of functional clusters in the *C. elegans* nervous system. Unsupervised Hebbian learning can be a self-optimization process to bring the initial network into a state of better generalization to new patterns.

Another important application of the study of complexity is the overarching field of **global issues** which is addressed in the following two papers. Nowhere does the importance of complexity science show as clearly as in issues of global ecological and economical systems which do not yield to simplistic treatments if one seeks them to be relevant. Amongst such topics, *sustainability* is of particular importance for the future of society and organized humanity. Molina-Perez et al. bring together an arsenal of methods to address the challenges of such models. An earlier model for the interplay of economical and ecological parameters, solved by constrained optimization, is considered under different policy regimes, uncertain stressors, and multiple experimentation with different elasticity parameters; it is subjected to machine-learning based clustering analysis for the parameters that produce stable vs. unstable regimes. All components, complex modeling, optimization, machine learning, and data mining work together to obtain a picture not only of how the system behaves on the whole, but also what policies should be enacted to increase the chance of obtaining desirable results.

Grisogono discusses how artificial intelligence (AI) could help tackle global complex problems, such as climate change, collapsing ecosystems, international conflicts, extremism, public policy, economics, and governance. These problems require

decision-making to attempt solutions or improvements. AI has the potential of mitigating failures and limitations in human decision-making, leading to a balance between robustness and adaptability. Nevertheless, there are also risks and drawbacks in AI, so its proper use and development should be discussed in detail.

Our research topic also shows how the concepts of complexity and self-organization are helpful in studying biological systems. Two above-mentioned papers Morales and Froese and Ramos-Fernandez et al. have contributions to both cognitive science and biology.

Our understanding of morphogenesis, as an important process of natural development, has advanced recently in part because of computing power, data availability, and algorithms. Thus, Pastor-Escuredo and del Álamo explore how computation is contributing to developmental biology. Computational models and simulations are proving useful to unravel the dynamic and multi-scale nature of morphogenesis. Machine learning and in particular deep learning architectures are promising in this respect. There are also potential applications for tissue engineering, identification of therapeutic targets, and synthetic life.

Finally, Casiraghi and Schweitzer address questions related to **computational social science**. In particular, they propose a method for improving the robustness of online social networks. Their aim is to prevent drop-out cascades of users. This is done using strategies to influence particular agents, reducing their probability of leaving the network, and thus considerably reducing drop-out cascades and increasing robustness.

As topics, complexity and self-organization have worked their way out of an exotic niche into the center of human activities. In contrast to the traditional reductionist treatment of scientific investigations, in today's crosstalk of disciplines, there is no field of human endeavor or study that can be considered in isolation. Understanding complexity has become a crucial skill in studying how the interacting levels of organismic function, society, ecological, and economical webs lead to a functioning whole—or to its disintegration. The richness of the contributions to this research topic serves as a showcase for the width and variety of tools and viewpoints that are being marshaled to this purpose.

AUTHOR CONTRIBUTIONS

All authors listed have made a substantial, direct and intellectual contribution to the work, and approved it for publication.

Conflict of Interest: The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2021 Gershenson, Polani and Martius. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.



Dynamical Inference of Simple Heteroclinic Networks

Maximilian Voit and Hildegard Meyer-Ortmanns*

Department of Physics and Earth Sciences, Jacobs University Bremen, Bremen, Germany

Heteroclinic networks are structures in phase space that consist of multiple saddle fixed points as nodes, connected by heteroclinic orbits as edges. They provide a promising candidate attractor to generate reproducible sequential series of metastable states. While from an engineering point of view it is known how to construct heteroclinic networks to achieve certain dynamics, a data based approach for the inference of heteroclinic dynamics is still missing. Here, we present a method by which a template system dynamically learns to mimic an input sequence of metastable states. To this end, the template is unidirectionally, linearly coupled to the input in a master-slave fashion, so that it is forced to follow the same sequence. Simultaneously, its eigenvalues are adapted to minimize the difference of template dynamics and input sequence. Hence, after the learning procedure, the trained template constitutes a model with dynamics that are most similar to the training data. We demonstrate the performance of this method at various examples, including dynamics that differ from the template, as well as a regular and a random heteroclinic network. In all cases the topology of the heteroclinic network is recovered precisely, as are most eigenvalues. Our approach may thus be applied to infer the topology and the connection strength of a heteroclinic network from data in a dynamical fashion. Moreover, it may serve as a model for learning in systems of winnerless competition.

Keywords: inference, heteroclinic networks, learning, metastable states, winnerless competition

OPEN ACCESS

Edited by:

Carlos Gershenson,
National Autonomous University of
Mexico, Mexico

Reviewed by:

Guoyong Yuan,
Hebei Normal University, China
Joseph T. Lizier,
University of Sydney, Australia

*Correspondence:

Hildegard Meyer-Ortmanns
h.ortmanns@jacobs-university.de

Specialty section:

This article was submitted to
Dynamical Systems,
a section of the journal
Frontiers in Applied Mathematics and
Statistics

Received: 19 September 2019

Accepted: 25 November 2019

Published: 10 December 2019

Citation:

Voit M and Meyer-Ortmanns H (2019)
Dynamical Inference of Simple
Heteroclinic Networks.
Front. Appl. Math. Stat. 5:63.
doi: 10.3389/fams.2019.00063

1. INTRODUCTION

When the unstable manifold of a saddle fixed point intersects the stable manifold of another saddle this is called a heteroclinic orbit. A heteroclinic network is a set of multiple saddles that are connected this way. Non-linear dynamics of heteroclinic networks are frequently found in ordinary differential equations under certain constraints like symmetries [1] or delay [2]. They are predicted in models of coupled phase oscillators [2, 3], vector models [2], pulse-coupled oscillators [4] and models of winnerless competition (WLC) [5]. Applications are manifold and range from social [6, 7] and ecological [5, 8] systems, to computation [4] and neuronal [9–13] networks. In particular, heteroclinic sequences in models of winnerless competition predict transient dynamics that share features with cognitive dynamics [5, 9–11, 13, 14]. Cognitive dynamics, or more generally, brain dynamics proceeds via sequential segmentation of information that is manifest in sequences of encephalography (EEG)-microstates [15] which are brief periods of stable scalp topography with a quasi-stationary configuration of the scalp potential field. Transitions between EEG-microstates have been modeled by epsilon-automata [16], for example, or by noisy network attractor models [17], of which the latter are closely related to heteroclinic networks. Such sequences of metastable states are observed on different time scales, ranging from milliseconds

to seconds [18]. In addition, these sequences may be nested as reflected in so-called chunking dynamics [10] when, for example, slow oscillations of neuronal activity modulate fast oscillations modulate even faster ones. On a formal level, the “events” in WLC are described as saddle equilibria (with one or higher-dimensional unstable manifolds), connected via heteroclinic orbits which facilitate transitions among the saddles [1]. The orbits can form heteroclinic sequences, cycles, or even whole networks with saddles as nodes and heteroclinic connections as edges. Specifically, heteroclinic networks are considered in this paper.

In an abstract representation, such sequences of metastable states can be modeled as a sequence of symbols, each representing one state with discrete state-transitions between them as in finite-state machines. In contrast, heteroclinic dynamics captures both, the sequence of states and autonomous smooth transitions between them as they exist in a physical realization.

How to construct a heteroclinic network with a certain topology has been well-studied, e.g., in references [19–21]. Moreover, with the perspective of engineering oscillators as noise driven heteroclinic cycles, the influence of different parameters on the dynamics has been investigated in Horchler et al. [22]. This way, versatile generators of repetitive patterns may be constructed by designing suitable heteroclinic networks.

In this article we address the inverse problem: Given the time series of a dynamics that was generated by a heteroclinic network, we propose how to infer the topology and the eigenvalues of this network. Related studies have been conducted for example by Selskii and Makarov [23]. The authors focus on how a learning process synchronizes the dynamics of heteroclinic cycles by adapting the expanding eigenvalues only. In Calvo Tapia et al. [24], this approach was extended by an additional step that identifies the sequence of saddles in a discrete manner, but it is still limited to circular topologies. With focus on the sequential memory in neural systems, Seliger et al. [25] proposed a learning mechanism for sequences of images based on winnerless competition. In their model, the learning mechanism that alters and adapts the competition matrix is realized via delay differential equations.

In this paper, we present a method that infers the topology and all eigenvalues of a so-called *simple* [26] heteroclinic network from time series data, generated by a heteroclinic network. Note that “simple” here does not refer to the topology, but to the type of heteroclinic network: Heteroclinic orbits of simple heteroclinic networks are contained in two-dimensional fixed-point subspaces, so that (for a suitably chosen coordinate system) the saddles lie on the coordinate axes. If the input was generated by a heteroclinic network, the time series of the process switches between metastable states, which manifest themselves in the data as accumulation points if the sampling rate was constant. Otherwise, if for a given time series the generating dynamics is not known, but the series shows such features of metastable states, the generation via a heteroclinic network would be a first conjecture. The inference is realized as a continuous dynamical process that alters the parameters of a template system. At the end of the process, this template system generates the same sequence of metastable states that was presented to it. The

method may thus be considered from various perspectives: As a data analysis/inference tool, as a tool for engineering purposes, and as a model of a learning process in the context of winnerless competition.

The remainder of this article is structured as follows. In section 2, we describe the method by introducing the template system and defining the learning dynamics. Additionally, we give a first demonstration of the method at a simple example, the Kirk-Silber network. Subsequently, we present increasingly complex networks in section 3 to highlight different aspects and possible obstacles in the application of the method. We conclude in section 4.

2. THE LEARNING DYNAMICS

Suppose we have an input signal $\mathbf{y}(t) \in \mathbb{R}^N$ that was generated by a simple heteroclinic network. In this case, the multidimensional time series has accumulation points (representing the metastable states) that lie on the coordinate axes in the positive hyperoctant (if necessary, after a suitable rotation). Moreover, we assume normalization, so that these accumulation points are essentially the unit vectors \mathbf{e}_i for $i \in \{1, \dots, N\}$. Our goal is to construct a system (described by ODEs) which generates a signal that resembles this input. To this end, we employ the idea that the ODEs of a simple heteroclinic network have a certain form as described in section 2.1. To mimic the dynamics of the input for a specific system, these ODEs of Equation (1) below merely have to be adjusted in their parameters. We call this adjustment (the incremental changes of the eigenvalues) the learning dynamics, defined in section 2.2. Afterwards, we demonstrate this method at a simple example, the Kirk-Silber network.

2.1. Template System

In the following, we describe the template system, which after training should reproduce the input sequence. Consider an input sequence $\mathbf{y}(t) \in \mathbb{R}^N$ with N accumulation points (representing the metastable states) on the coordinate axes in the positive hyperoctant (Depending on the context, the variable y may represent, for example, species concentrations, cognitive information, or whatever physical meaning the temporary winner in this case of WLC has). To produce such a sequence by a simple heteroclinic network, N dimensions are required, as saddle fixed points are located only on the coordinate axes. We thus propose as template

$$\begin{aligned} \dot{x}_i = x_i & \left(-a_{i,i} \left(\frac{1}{2} + \frac{x_i^2}{b_i^2} \right) + \sum_j (a_{j,i} + \frac{a_{i,i}}{2}) \frac{x_j^2}{b_j^2} \right) \\ & + \sigma \eta_i(t) \quad \forall i \in \{1, \dots, N\}, \end{aligned} \tag{1}$$

where $\mathbf{x} \in \mathbb{R}^N$ describes the state, and $a_{i,j} \in \mathbb{R}$ and $b_i \in \mathbb{R}, b_i > 0$ are free parameters that will be subject to learning. Indices i, j, k are always assumed $\in \{1, \dots, N\}$. The parameter $\sigma \geq 0$ determines the noise strength, and $\eta_i(t)$ is white noise with zero mean and unit variance. This system has N equilibria $\xi_i = \{\mathbf{x} \in \mathbb{R}^N : x_i = b_i > 0, x_j = 0 \forall j \neq i\}$ with only a single item $x_i > 0$ active. Moreover, the eigenvalues of the Jacobian of Equation (1)

evaluated at ξ_i are $a_{i,j}$, and the corresponding eigenvectors are $e_j = \{x \in \mathbb{R}^N : x_j = 1, x_k = 0 \forall k \neq j\}$. Precisely these properties are the reason for choosing the very form of Equation (1). It is the lowest order realization that has the aforementioned properties and \mathbb{Z}_2^N reflection symmetry. A second order realization would in principle be also possible, but lacks this symmetry and may show divergent dynamics as soon as $x_i < 0$ for some i . Equation (1) may be understood as a generalization of the ODEs of the simplex method [19] that makes all eigenvalues and the saddle positions directly available as parameters. These ODEs can be retrieved from Equation (1) by setting $b_i = 1$ and $a_{i,i} = -2$ for all i .

Note that Equation (1) is equivariant under reflection symmetry \mathbb{Z}_2^N . As a result, the coordinate planes $P_{i,j} = \{x \in \mathbb{R}^N : x_k = 0 \forall k \notin \{i,j\}\}$ are invariant sets. Thus, when the eigenvalues of two equilibria ξ_i and ξ_j fulfill $a_{i,j} < 0, a_{j,i} > 0, a_{i,i} < 0$, and $a_{j,j} < 0$, there exists a heteroclinic orbit from ξ_i to ξ_j within $P_{i,j}$ [27, 28]. Furthermore, the hyperoctants (e.g., \mathbb{R}_+^N) are invariant sets, so in the following we assume w.l.o.g. all components of x (and y) to be positive at all times. For simplicity, we also assume that the input sequence $y(t)$ is normalized so that the accumulation points are the unit vectors $e_i \in \mathbb{R}^N$. We thus fix $b_i = 1 \forall i$.

2.2. Definition of the Learning Dynamics

The principal idea is to make the template system follow the input signal $y(t) \in \mathbb{R}^N$ by coupling it linearly into the template. If we know also the system that generates the input, the learning may be performed online, so that the signal is learned while it is generated. When the generating system is known, the setup can be seen as a master-slave coupling, as there is no coupling back from the template to the generating system. Thus, Equation (1) is extended by a coupling term

$$d_t x_i = x_i (\dots) + \sigma \eta_i(t) + \vartheta (y_i - x_i) \tag{2}$$

with y_i the i th component of the input signal and ϑ the strength of the coupling. Empirically, for ϑ large enough the coupling draws the template dynamics close to the input as desired, even under the influence of noise. For mutual coupling of two identical heteroclinic networks without noise such an effect may be anticipated via the master stability function approach [29]: The mode corresponding to the synchronized manifold has the original eigenvalues, say λ_j ; the transverse mode has eigenvalues $\lambda_l - 2\vartheta$, and perturbations away from the synchronized manifold thus decay if the coupling ϑ is large enough. A detailed discussion of this synchronization (in the sense that $\|x(t) - y(t)\| \rightarrow 0$ for $t \rightarrow \infty$) for linearly coupled heteroclinic cycles will be given elsewhere (Voit and Meyer-Ortmanns, in preparation).

Even with coupling, however, small differences between master and slave remain as long as the two systems are not identical. The key point is therefore the following: When the trajectory is in the vicinity of saddle ξ_i , it is the N eigenvalues $a_{i,j}$ at ξ_i that determine the time evolution of the concentrations x_j near ξ_i . If there is a difference $(y_j - x_j) > 0 (< 0)$ while the systems are close to ξ_i , it is therefore the eigenvalue $a_{i,j}$ that has to

be increased (decreased) to match the eigenvalue underlying the signal. This is realized by the learning rule

$$d_t a_{i,j} = (1 - \delta_{ij}(1 + \rho)) \gamma \vartheta (y_j - x_j) \exp\left(-\left(b_i - \frac{y_i + x_i}{2}\right)^2 \zeta b_i\right), \tag{3}$$

where $\gamma > 0$ is the learning rate, δ_{ij} the Kronecker delta. The first terms are precisely the scaled dependence on the deviation of the current dynamics x from the original y . By taking along a factor of ϑ , the learning rate γ becomes independent of the coupling strength. The exponential term is a Gaussian ensuring that the changes of eigenvalues are local to the saddle these eigenvalues are associated with: The difference $b_i - \frac{y_i + x_i}{2}$ becomes small precisely when the average of the dynamics of the input and the template is close to the location of the saddle. It should be noticed that here the structure of simple heteroclinic networks enters in that it suffices to measure the i th component only, since regularly never a situation occurs where two coordinates i, j simultaneously strongly differ from zero such that $x_i \approx b_i$ and $x_j \approx b_j$ at the same time. The range of this localization is adjusted by the parameter $\zeta > 0$. The b_i -dependence is kept in the exponent for cases where $b_i \neq 1$ to adjust the size of the neighborhood of the saddle. Note that the situation for the radial eigenvalues is different. It is necessary to use coordinates local to the saddle, which for non-radial components are just the global ones. The radial component, however, is transformed to $\tilde{x}_i = b_i - x_i$ (equivalently for y_i) in local coordinates, so that the sign of the learning rule has to be inverted. We therefore require $\rho > 0$ and usually will choose $\rho \geq 1$ since radial eigenvalues empirically converge slower than eigenvalues associated with the other directions.

2.3. Inferring a Single Eigenvalue

We proceed by illustrating the method introduced above with a Kirk-Silber network [30]. This is a simple heteroclinic network consisting of two heteroclinic cycles that share a common edge, c.f. **Figure 1**. Suppose that the ODEs of the master system $d_t y$ are known and of the form of Equation (1). The slave system (the template) naturally is Equation (2), and we assume to know all eigenvalues $a_{i,j}$ but $a_{2,3}$, which is different from its value in the master system $a_{2,3m}$. The effect of the linear coupling is to continuously counteract this difference, but ultimately the learning dynamics of Equation (3) leads to the convergence $a_{2,3} \rightarrow a_{2,3m}$, and the contribution of the coupling term in Equation (2) vanishes, c.f. **Figure 2**. Note that the learning takes place whenever the template system visits ξ_2 and x_3 differs from y_3 . During the remaining time, the differences between x and y are due to the differing noise realizations in both systems, which also makes both dynamics diverge as soon as the coupling is removed at $t = 1,500$. Afterwards the fact that the template (slave system) on its own has the same dynamics as the master system is clear from its value of $a_{2,3} = a_{2,3m}$ on the one hand, and the statistics (of visits to ξ_3 vs. ξ_4) on the other hand. It might be beneficial to delay the start of learning in order to allow initial transients to decay (this is not necessary when the initial condition is close to the heteroclinic network).

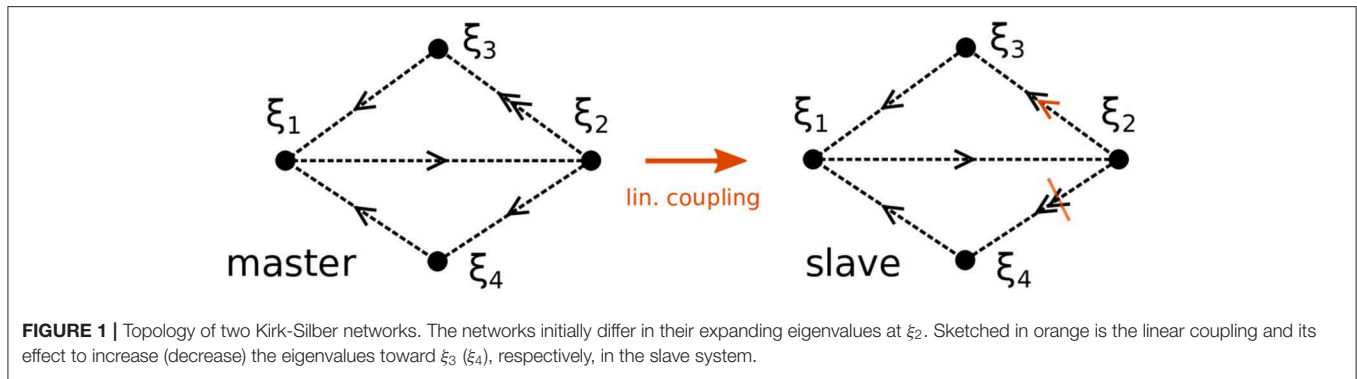


FIGURE 1 | Topology of two Kirk-Silber networks. The networks initially differ in their expanding eigenvalues at ξ_2 . Sketched in orange is the linear coupling and its effect to increase (decrease) the eigenvalues toward ξ_3 (ξ_4), respectively, in the slave system.

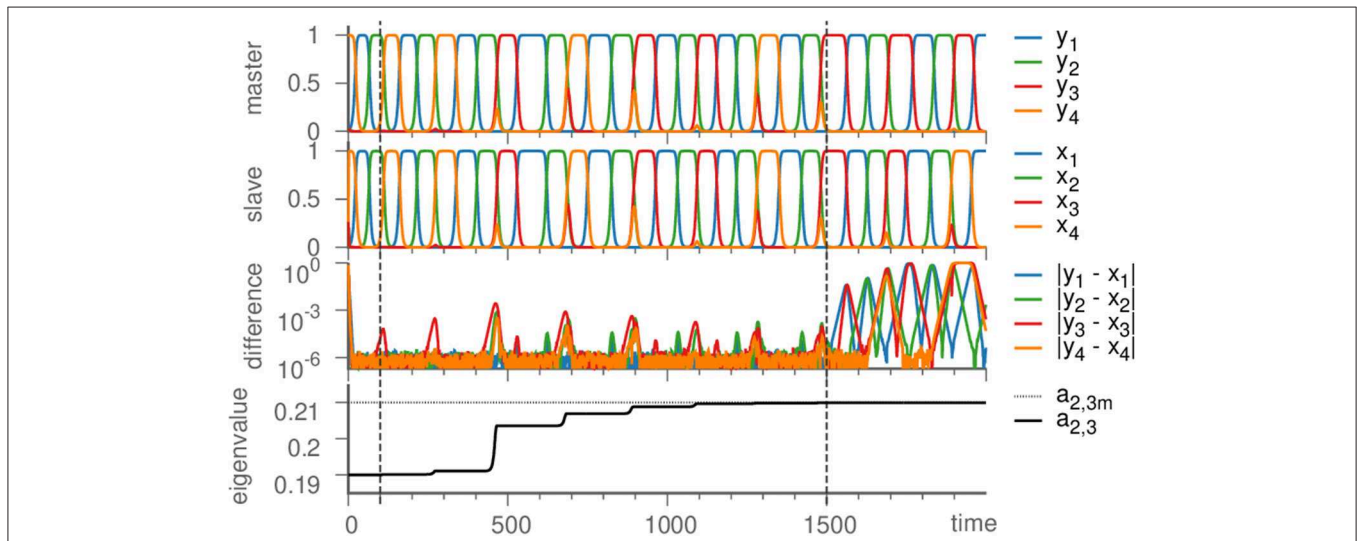


FIGURE 2 | Learning dynamics near the Kirk-Silber network. From top to bottom: Dynamics of the master system, the slave system, their difference, and the eigenvalue $a_{2,3}$ are plotted against time. Vertical dashed lines mark the beginning of learning ($t = 100$) and its end ($t = 1,500$), at which time also the coupling is turned off ($\rightarrow \vartheta = 0$), so that due to different noise realizations the system states slowly diverge. Parameters are $\vartheta = 1$, $\gamma = 0.5$, $\zeta = 50$, $\sigma = 10^{-6}$; eigenvalues are chosen as $a_{rad.} = -1$, $a_{contr.} = -0.22$, $a_{exp.} = 0.2$, and $a_{2,3m} = 0.21$.

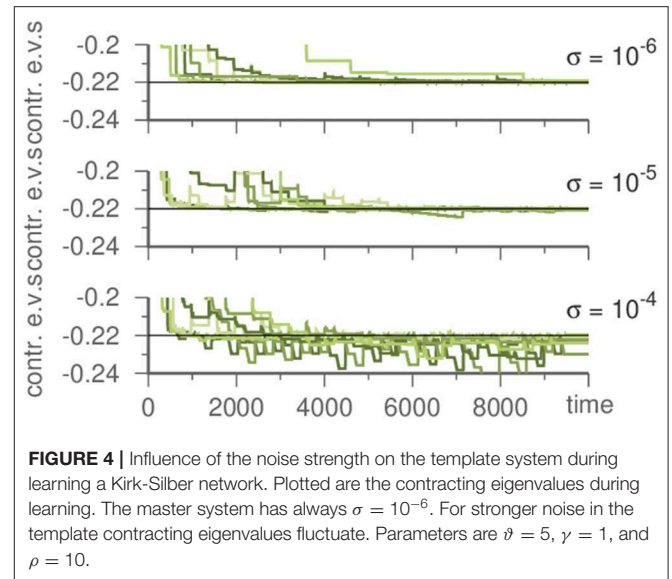
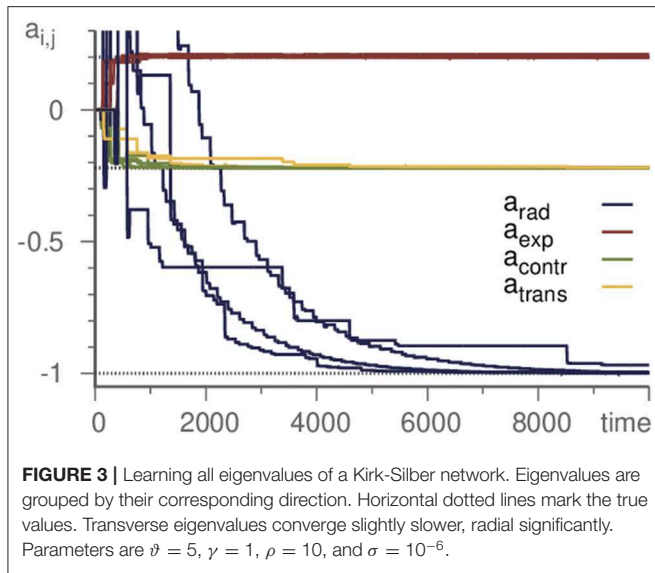
3. INFERRING HETEROCLINIC NETWORKS WITH INCREASING COMPLEXITY

In this section, we illustrate our method of the previous section by heteroclinic networks with increasingly complex features. As the simplest non-trivial topology we choose the Kirk-Silber network in section 3.1. At this example, we demonstrate how not only one, but all eigenvalues are recovered by the template without any prior knowledge. In addition, we point out how even the noise level may be captured by the learning method. Moreover, in section 3.2 we analyze the effect of a mismatch between the system that generated the input and the template system. Significant differences in the ODEs of the two systems strongly affect the convergence of radial eigenvalues, while the remaining eigenvalues are mostly inferred well. Furthermore, in section 3.3 we focus on larger networks with more complex topologies. We both probe our method at a highly regular, hierarchical heteroclinic network which exhibits two time scales, and construct a random heteroclinic network (by the simplex

method) to generate the input and reconstruct its topology by learning. The latter example underlines the role of noise in how extensively the heteroclinic network is explored, especially in the case of an irregular topology with heterogeneous preferences of heteroclinic connections.

3.1. Inferring All Eigenvalues and the Noise Level

As the basic example in section 2, we demonstrated the successful inference of a single eigenvalue of a Kirk-Silber network. Actually, however, all eigenvalues may be inferred at the same time. Thus it is possible to start with a template with unbiased randomly or uniformly chosen parameters and infer the whole topology of a simple heteroclinic network. As demonstration, again we choose the Kirk-Silber network and initialize all eigenvalues as 0. The learning method then infers the values of the generating system, c.f. **Figure 3**. It is convenient to distinguish the different kinds of eigenvalues and refer to them by standard terminology [31] according to their respective



eigendirection (radial, contracting, expanding, and transverse). Note that transverse eigenvalues converge only slightly slower than contracting and expanding ones, while radial eigenvalues pose a greater difficulty and converge much slower. Therefore, accelerating the learning process for radial eigenvalues (by choosing $\rho \geq 1$) may be useful to moderate this effect.

Up to this point we neglected a careful discussion of the influence of noise, although noise has to be present in the generating system to sustain the switching between saddles. If the template should reproduce this switching after the input signal is switched off, it must also be subject to noise. As the noise intensity influences the pace of switching, it ought to be the same in both systems.

To discover the original noise intensity we exploit its characteristics. If the noise is lower in the template than it was in the generating system, this has no noticeable effect. However, if it is stronger in the template, contracting eigenvalues fluctuate during learning, c.f. **Figure 4**. Thus, by performing multiple learning trials with decreasing σ in Equation (1), the correct noise level (as inherent in the input) may be recovered.

In summary, after decoupling the template from the master system with identical dynamics, the time evolution of both systems is exactly the same if there is no noise and identical initial conditions have been used. Otherwise, i.e., under the influence of noise and depending on the initial conditions, the statistics and sequence of visited saddles in both decoupled time evolutions remain the same, but the dynamics differs in details.

3.2. Mismatched Template

In the examples above the input signal is generated by a system of ODEs that has the same form as the template. Otherwise, if the input stems from a different implementation of a simple heteroclinic network, the question arises of how this mismatch between template and generating system impacts the inference. In the following we pursue this question, as it is crucial in view of

the fact that for a realistic inference task the form of the original ODEs is usually unknown.

Time continuous models of population dynamics are commonly derived as a mean-field approximation [32] of reaction equations that describe interactions at the level of individuals. One basic example of such a continuous model is the May-Leonard model [33]. It contains a heteroclinic cycle that is also known as the Busse-Heikes cycle [34], generated by the ODEs

$$d_i x_i = x_i(1 - x_i - b x_{i+1} - c x_{i+2}), \tag{4}$$

where $0 < c < 1 < b$, $b - 1 > 1 - c$, and $i \in \mathbb{Z}_3 = \{1, 2, 3\}$ cyclic. The variables x_i represent population densities, thus they are restricted to the positive octant $\mathbb{R}_{\geq 0}^3$. By a variable transformation ($x_i \rightarrow \sqrt{x_i}$) the Guckenheimer-Holmes cycle [35] emerges, which matches the form of the template. The original Busse-Heikes cycle, however, does not; it has second-order terms instead of third-order ones¹.

Nevertheless, our method is able to infer the Busse-Heikes cycle. With the default parameters ($\vartheta = 1$, $\gamma = 0.5$, $\zeta = 50$, and $\rho = 10$), we find that the radial eigenvalues fluctuate strongly, but the eigenvalues of the remaining directions converge approximately to their true values. Choosing a low value of ρ (e.g., $\rho = 0.2$) reduces the fluctuations of the radial eigenvalues. The resulting template follows a heteroclinic cycle with the same topology, but different shape of the approach toward the saddles, c.f. **Figure 5A**.

Since the heteroclinic cycles that we considered so far do not contain transverse directions and we want to analyze the effect of a mismatch also on the transverse eigenvalues, we modified a

¹Commonly, the terms “Busse-Heikes cycle” and “Guckenheimer-Holmes cycle” are used synonymously, as the heteroclinic cycles (as objects in phase space) are diffeomorphic to each other. In this article, however, we specifically distinguish the two different ODE systems by these terms.

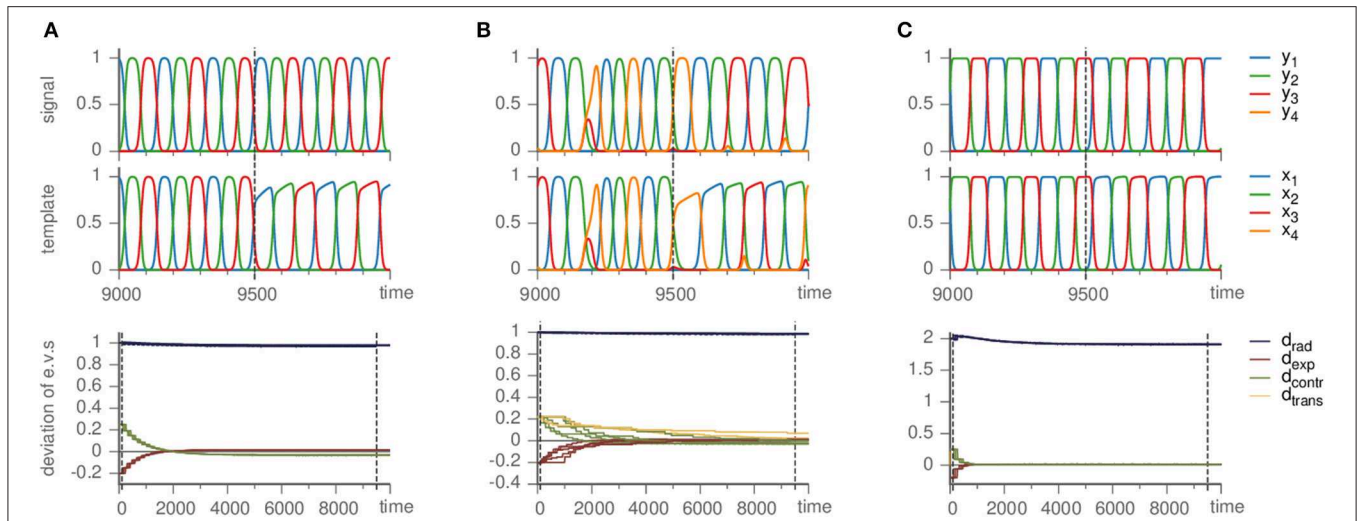


FIGURE 5 | Effects of mismatched templates. Upper panels show the dynamics of the signal and template against time. Lower panels display the deviation of the eigenvalue in the template from their value a_{ijm} in the master system, i.e., $d_{ij} = a_{ij} - a_{ijm}$, sorted by the kind of direction they correspond to (radial, expanding, contracting, and transverse). Parameters are $\vartheta = 1$, $\gamma = 0.5$. **(A)** Busse-Heikes cycle, with $b = 1.25$, $c = 0.8$, and $\rho = 0.2$. **(B)** Kirk-Silber network with only second-order terms, using $\rho = 0.1$. **(C)** Guckenheimer-Holmes cycle with higher order terms, using $a = 1$, $b = 1.25$, $c = 0.8$, $\alpha = 0.01$, $\beta = 0.01$, and $\rho = 1$.

Kirk-Silber system in the same way (so that its terms are second-order instead of third-order). The situation is quite similar; small ρ remedies the fluctuation of radial eigenvalues, whereas the transverse eigenvalues converge toward their true values, only slower than expanding and contracting ones, c.f. **Figure 5B**.

One further possible mismatch of the template and the generating system is due to higher order terms. To check the effect, we modified a Guckenheimer-Holmes cycle by adding two fourth-order terms:

$$d_t x_i = x_i (1 - a x_i^2 - b x_{i+1}^2 - c x_{i+2}^2 - \alpha x_i^4 - \beta x_i^2 x_{i+1}^2). \quad (5)$$

The additional terms do not break the \mathbb{Z}_2^3 equivariance. Thus, for $|\alpha|$ and $|\beta|$ small, the heteroclinic cycle persists (it is structurally stable). While the second term affects the dynamics far from the saddles (for $\beta \neq 0$), the first one acts in their vicinity. More precisely, $\alpha \neq 0$ changes the position of the saddles and also the eigenvalues, so $|\alpha| \ll 1$ is necessary to maintain the heteroclinic cycle.

As long as the cycle persists, our method correctly identifies it and approximately infers the eigenvalues (independently on whether they are original or changed due to $\alpha \neq 0$) of the generating system, c.f. **Figure 5C**. Here, as in the other cases of mismatched templates, the eigenvalues corresponding to radial directions fluctuate and converge to values different from the ones in the generating system. More precisely, we observed radial eigenvalues to be only slightly negative, even though these directions should be definitely stable. In contrast to section 3.1 setting $\rho \geq 1$ is not helpful, but intensifies the problem. Instead, a possible remedy is to ensure that radial directions are stable by fixing the radial eigenvalues to a negative value (e.g., -1) from the beginning and keeping them at this value rather than changing them by the learning dynamics (by setting $\rho = 0$), see the following example.

3.3. Inferring Larger Regular and Irregular Networks

The example networks presented up to this point were rather simple in their topology, involving four saddles at most. Larger networks may pose additional challenges for inference, as we point out by the following two examples: one is a highly regular, hierarchical heteroclinic network with nine nodes; the other one is a random heteroclinic network composed of 12 nodes with heterogeneous in- and out-degrees.

In Voit and Meyer-Ortmanns [36], we constructed a heteroclinic network \mathcal{H} that is hierarchically structured. It consists of three small heteroclinic cycles (SHCs) that constitute the saddles of a large heteroclinic cycle (LHC). This hierarchy is produced by a difference of the expanding eigenvalues associated with connections belonging to one SHC vs. connections between different SHCs. The structural hierarchy translates to a hierarchy in time scales, which amounts to the modulation of fast oscillations by slower ones. The network \mathcal{H} obeys a $\mathbb{Z}_3 \times \mathbb{Z}_3$ symmetry. Thus it is highly regular, as is its dynamics. All saddles are visited equally often, and all SHCs dominate with the same frequency.

We apply our inference method to the dynamics generated by the very system described in Voit and Meyer-Ortmanns [36], c.f. **Figure 6**. It thus deviates from the template dynamics by containing only second-order terms compared to the third-order terms of the template. This mismatch leads to a deviation of the inferred radial eigenvalues from the real ones, c.f. **Figure 6D**, just as expected from the previous section. Nevertheless, the topology of \mathcal{H} and its structural hierarchy (manifest as the difference between the two kinds of expanding eigenvalues in the small and large heteroclinic cycles) is inferred correctly. The resulting dynamics of the template thus reproduces the same sequence of saddles visited as the

original system. Frequency and amplitude, however, are not recovered accurately.

As an alternative, we fixed the radial eigenvalues to their true value $a_{i,i} = -1$ and subjected only $a_{i,j}$ for $i \neq j$ to the learning method (setting $\rho = 0$). For these entries, the error of the inference is comparable to the previous situation, c.f. **Figure 6E**. In addition, due to employing the true radial eigenvalues, the resulting dynamics does not only reproduce the sequence of the visited saddles correctly, but also recovers the frequency and amplitude of the oscillations to a good agreement, c.f. **Figure 6H**.

For the random heteroclinic network, we generated a 12 node Erdős-Rényi graph with edge probability 0.2, without self-loops. Subsequently, two-loops were removed by deleting one of the edges each, whilst ensuring that the in- and out-degree at all nodes is at least one. **Figure 7** depicts the topology of the resulting graph \mathcal{G} .

From the graph we generated the heteroclinic network by employing the same form of ODEs as in the template Equation (1) and choosing the eigenvalues $a_{i,j}$ from the adjacency matrix A in the following way:

$$a_{i,j} = \begin{cases} X \in (0.4, 0.6) & \text{for } A_{i,j} = 1 \\ X \in (-1.1, -0.9) & \text{for } A_{i,j} = 0 \end{cases} \quad (6)$$

with X a random variable taken uniformly from the specified interval. The choice of these intervals is arbitrary to some extent. Mainly, eigenvalues in expanding directions ($A_{i,j} = 1$) must be positive, while contracting, radial and transverse eigenvalues must be negative. Furthermore, for the heteroclinic network as

a whole to be attractive, “contraction must surpass expansion”. The size of the intervals controls the degree of heterogeneity in the preference of heteroclinic orbits. Overall, this process is thus an adapted version of the simplex method [19], which describes how to construct a simple heteroclinic network for a given graph.

For our choice of intervals, the expanding eigenvalues differ sufficiently, so that the system dwells more frequently in some parts of the network than in others. The relevance of this becomes especially clear once the learning method is applied. Strong noise is required to infer all parts of the heteroclinic network, c.f. **Figure 8**. Then, however, also the inferred eigenvalues fluctuate

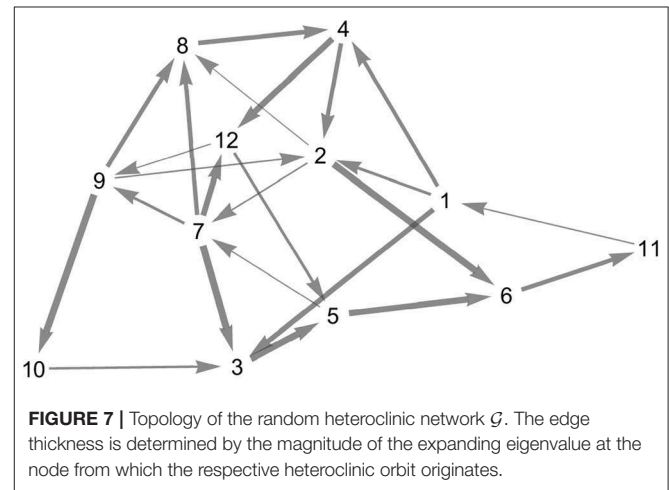


FIGURE 7 | Topology of the random heteroclinic network \mathcal{G} . The edge thickness is determined by the magnitude of the expanding eigenvalue at the node from which the respective heteroclinic orbit originates.

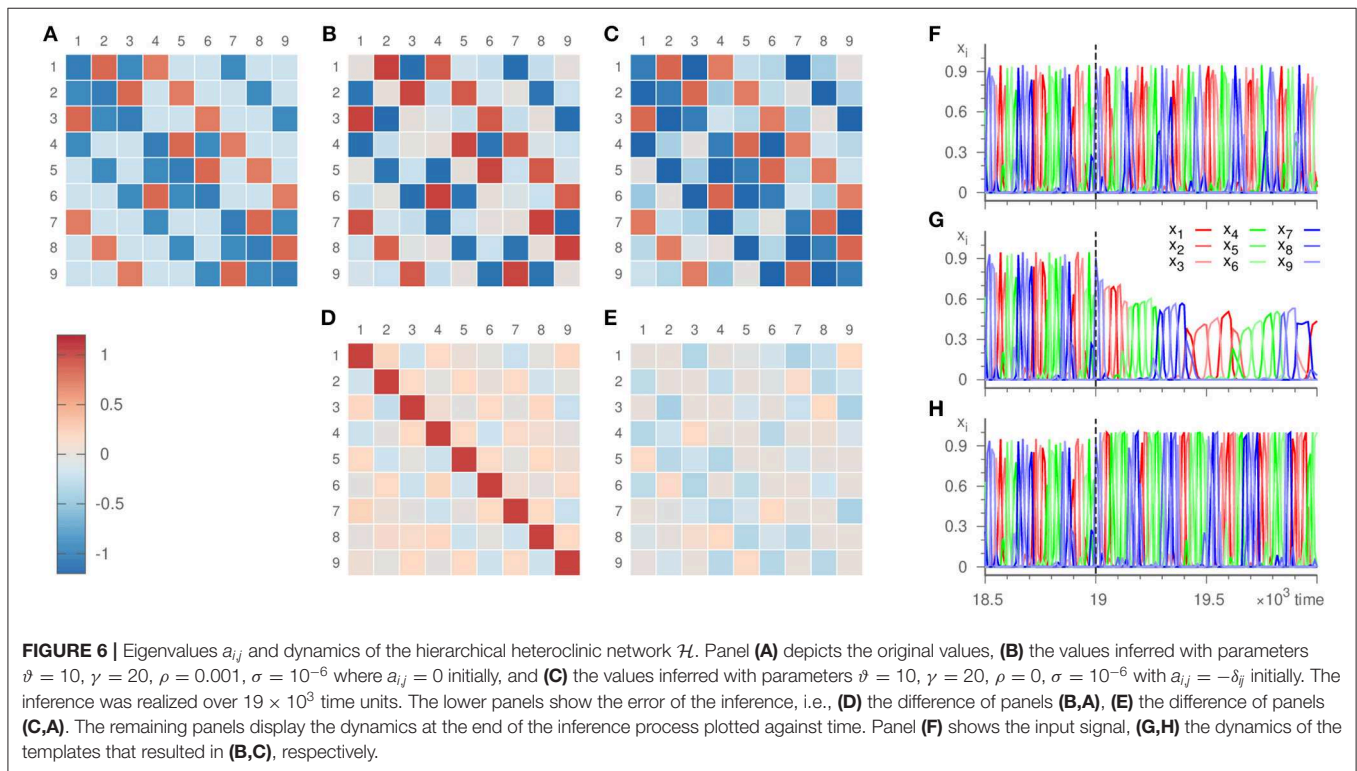
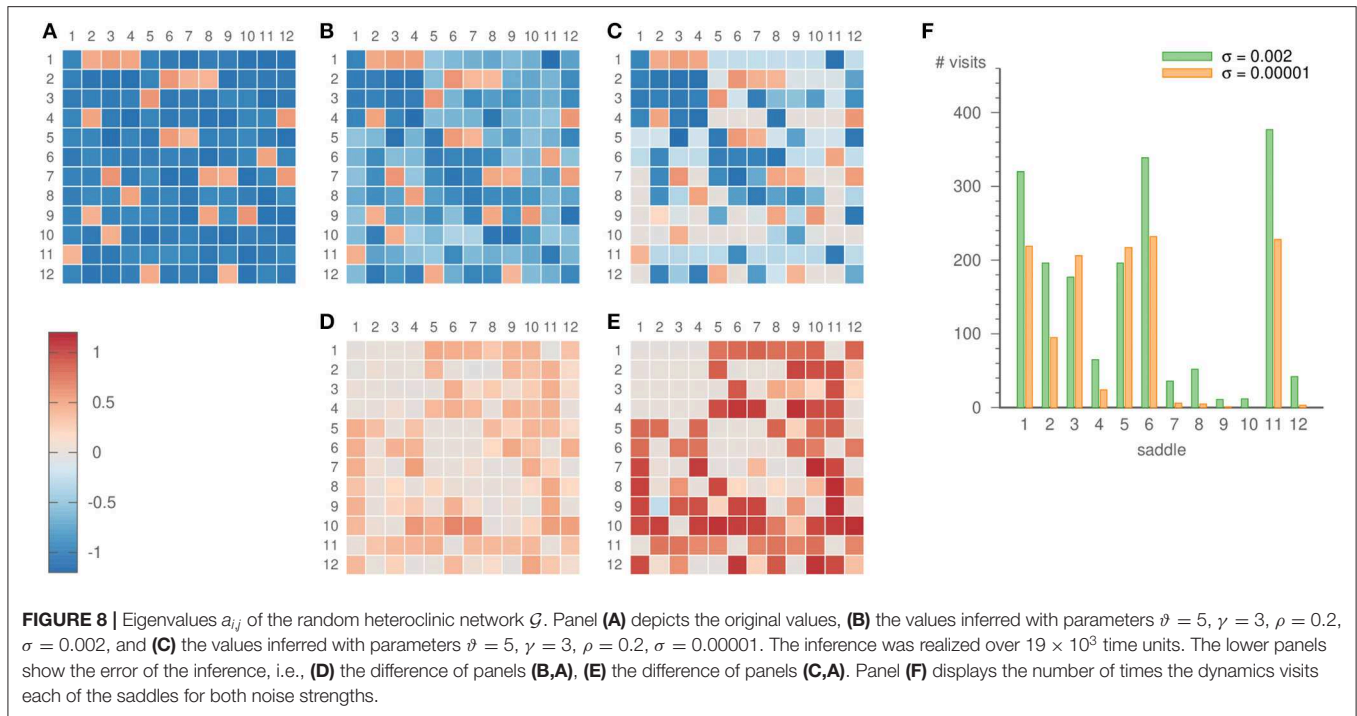


FIGURE 6 | Eigenvalues $a_{i,j}$ and dynamics of the hierarchical heteroclinic network \mathcal{H} . Panel **(A)** depicts the original values, **(B)** the values inferred with parameters $\vartheta = 10, \gamma = 20, \rho = 0.001, \sigma = 10^{-6}$ where $a_{i,j} = 0$ initially, and **(C)** the values inferred with parameters $\vartheta = 10, \gamma = 20, \rho = 0, \sigma = 10^{-6}$ with $a_{i,j} = -\delta_{ij}$ initially. The inference was realized over 19×10^3 time units. The lower panels show the error of the inference, i.e., **(D)** the difference of panels **(B,A)**, **(E)** the difference of panels **(C,A)**. The remaining panels display the dynamics at the end of the inference process plotted against time. Panel **(F)** shows the input signal, **(G,H)** the dynamics of the templates that resulted in **(B,C)**, respectively.



strongly. For low noise levels, on the other hand, some saddles are visited only rarely. For example, saddle 10 is not visited at all for $\sigma = 10^{-5}$, but 12 times for $\sigma = 0.002$, c.f. **Figure 8F**. This scarcity of visits to certain saddles is one factor that may lead to a comparatively bad inference of their eigenvalues, c.f. **Figure 8E**. However, other factors such as the topology of the heteroclinic network and the recent history of the trajectory before arriving at a certain saddle play a role as well.

For practical applications the obvious trade-off between exploring the whole network (high noise level) and the quality of the inferred eigenvalues (low noise level) is of minor importance. Indeed, for weak noise it would be impossible to infer some of the saddles, but the actual dynamics neither visits these saddles.

Besides the noise strength, the length of the input signal needs to be taken into account. Longer input is beneficial, as weakly attached parts of the networks get visited more often. If the input is too short, only the most probable cycles of the heteroclinic network become inferred. For example, running the inference for merely 1500 time units, we observed the resulting heteroclinic network settle to the cycles $1 \rightarrow 2 \rightarrow 6 \rightarrow 11 \rightarrow 1$, or $1 \rightarrow 3 \rightarrow 5 \rightarrow 6 \rightarrow 11 \rightarrow 1$.

4. DISCUSSION

In summary, we have introduced a novel method of learning simple heteroclinic networks. It is based on an unbiased template system of a heteroclinic network in combination with a learning dynamics that progressively alters the eigenvalues at the saddles. The system thereby dynamically infers the eigenvalues at all saddles and thus reconstructs the topology of the heteroclinic network that generated the signal. A key ingredient is the linear

coupling to the input signal, which forces the dynamics of the input onto the template. Only this enables the learning, which primarily takes place when the system visits the saddle equilibria. The trained template then reproduces sequences of metastable states most similar to the input time series.

We worked out the performance of this method for various examples, inferring all eigenvalues even in comparatively large heteroclinic networks. Moreover, we illustrated possible difficulties that the noise level or a mismatch of template and generating system can pose, for example. We pointed out strategies to handle them. A subtle point will be to achieve a deeper understanding of what determines the speed of learning the eigenvalues of saddles, that is, its dependence on the topology of the heteroclinic network, the noise level and other factors.

In view of engineering underlying heteroclinic networks from a given data set, our method provides a continuous counterpart to designing simple finite-state machines from given example data, as it automatically interpolates between subsequent maxima. If data of sequential switching between different metastable states suggest games of winnerless competition behind their generation, it would be natural to attempt a learning of rates at a first place (say in generalized Lotka-Volterra models), rather than a learning of eigenvalues. In simple heteroclinic networks it is the local information stored in the eigenvalues of the saddles that is sufficient to control and learn the time evolution of the dynamics in a desired way, bridging the global (non-local) distance between the different saddles. Therefore, as long as the assumed heteroclinic network is simple, one would learn the rates as a function of the learned eigenvalues, while the eigenvalues at the saddles are expressed in terms of the rates.

Simple heteroclinic networks are specific in the sense that the saddles lie on the coordinate axes, the phase space has a dimension that is given by the number of saddles, and together with the imposed symmetry one knows from the local information of the eigenvalues at one saddle at which saddle one ends up next. It is therefore sufficient to learn the eigenvalues (and thus mimic the local dynamics) in order to reproduce the global dynamics. In general (and in particular in the context of heteroclinic computing), the heteroclinic networks are non-simple and the dimension of phase space is lower than the number of saddles. It is an interesting open challenge to derive rules for learning non-simple heteroclinic networks and possibly combine these with the concept of heteroclinic computing.

DATA AVAILABILITY STATEMENT

All datasets generated for this study are included in the article/supplementary material.

REFERENCES

- Krupa M. Robust Heteroclinic Cycles. *J Nonlinear Sci.* (1997)7:129–76. doi: 10.1007/BF02677976
- Kori H, Kuramoto Y. Slow switching in globally coupled oscillators: robustness and occurrence through delayed coupling. *Phys Rev E.* (2001) 63:046214. doi: 10.1103/PhysRevE.63.046214
- Wordsworth J, Ashwin P. Spatiotemporal coding of inputs for a system of globally coupled phase oscillators. *Phys Rev E.* (2008) 78:066203. doi: 10.1103/PhysRevE.78.066203
- Schittler Neves F, Timme M. Computation by switching in complex networks of states. *Phys Rev Lett.* (2012) 109:018701. doi: 10.1103/PhysRevLett.109.018701
- Afraimovich V, Tristan I, Huerta R, Rabinovich MI. Winnerless competition principle and prediction of the transient dynamics in a Lotka–Volterra model. *Chaos Interdiscip J Nonlinear Sci.* (2008) 18:043103. doi: 10.1063/1.2991108
- Hauert C, Szabó G. Game theory and physics. *Am J Phys.* (2005) 73:405–14. doi: 10.1119/1.1848514
- Szabó G, Vukov J, Szolnoki A. Phase diagrams for an evolutionary prisoner's dilemma game on two-dimensional lattices. *Phys Rev E.* (2005) 72:047107. doi: 10.1103/PhysRevE.72.047107
- Nowak MA, Sigmund K. Biodiversity: bacterial game dynamics. *Nature.* (2002) 418:138–9. doi: 10.1038/418138a
- Rabinovich MI, Afraimovich VS, Varona P. Heteroclinic binding. *Dyn Syst.* (2010) 25:433–42. doi: 10.1080/14689367.2010.515396
- Rabinovich MI, Varona P, Tristan I, Afraimovich VS. Chunking dynamics: heteroclinics in mind. *Front Comput Neurosci.* (2014) 8:22. doi: 10.3389/fncom.2014.00022
- Rabinovich MI, Huerta R, Varona P, Afraimovich VS. Transient cognitive dynamics, metastability, and decision making. *PLoS Comput Biol.* (2008) 4:e1000072. doi: 10.1371/journal.pcbi.1000072
- Komarov MA, Osipov GV, Suykens JAK. Metastable states and transient activity in ensembles of excitatory and inhibitory elements. *EPL Europhys Lett.* (2010) 91:20006. doi: 10.1209/0295-5075/91/20006
- Fingelkurts A, Fingelkurts A. Information flow in the brain: ordered sequences of metastable states. *Information.* (2017) 8:22. doi: 10.3390/info8010022
- Afraimovich VS, Rabinovich MI, Varona P. Heteroclinic contours in neural ensembles and the winnerless competition principle. *Int J Bifurc Chaos.* (2004) 14:1195–208. doi: 10.1142/S0218127404009806
- Michel CM, Koenig T. EEG microstates as a tool for studying the temporal dynamics of whole-brain neuronal networks: a review. *NeuroImage.* (2018) 180:577–93. doi: 10.1016/j.neuroimage.2017.11.062
- Nehaniv CL, Antonova E. Simulating and reconstructing neurodynamics with Epsilon-automata applied to electroencephalography (EEG) microstate sequences. In: *2017 IEEE Symposium Series on Computational Intelligence (SSCI)*. Honolulu, HI (2017). p. 1–9. doi: 10.1109/SSCI.2017.8285438
- Creaser J, Ashwin P, Postlethwaite C, Britz J. Noisy network attractor models for transitions between EEG microstates. *arXiv:1903.05590* (2019).
- Van De Ville D, Britz J, Michel CM. EEG microstate sequences in healthy humans at rest reveal scale-free dynamics. *Proc Natl Acad Sci USA.* (2010) 107:18179–84. doi: 10.1073/pnas.1007841107
- Ashwin P, Postlethwaite C. On designing heteroclinic networks from graphs. *Phys Nonlinear Phenom.* (2013) 265:26–39. doi: 10.1016/j.physd.2013.09.006
- Field MJ. Heteroclinic networks in homogeneous and heterogeneous identical cell systems. *J Nonlinear Sci.* (2015) 25:779–813. doi: 10.1007/s00332-015-9241-1
- Ashwin P, Postlethwaite C. Designing heteroclinic and excitable networks in phase space using two populations of coupled cells. *J Nonlinear Sci.* (2016) 26:345–64. doi: 10.1007/s00332-015-9277-2
- Horchler AD, Daltorio KA, Chiel HJ, Quinn RD. Designing responsive pattern generators: stable heteroclinic channel cycles for modeling and control. *Bioinspir Biomim.* (2015) 10:026001. doi: 10.1088/1748-3190/10/2/026001
- Selskii A, Makarov VA. Synchronization of heteroclinic circuits through learning in coupled neural networks. *Regul Chaot Dyn.* (2016) 21:97–106. doi: 10.1134/S1560354716010056
- Calvo Tapia C, Tyukin IY, Makarov VA. Fast social-like learning of complex behaviors based on motor motifs. *Phys Rev E.* (2018) 97:052308. doi: 10.1103/PhysRevE.97.052308
- Seliger P, Tsimring LS, Rabinovich MI. Dynamics-based sequential memory: winnerless competition of patterns. *Phys Rev E.* (2003) 67:011905. doi: 10.1103/PhysRevE.67.011905
- Krupa M, Melbourne I. Asymptotic stability of heteroclinic cycles in systems with symmetry. II. *Proc R Soc Edinb Sect Math.* (2004) 134:1177–97. doi: 10.1017/S0308210500003693
- Hofbauer J, Sigmund K. *Evolutionary Games and Population Dynamics*. Cambridge: Cambridge University Press (1998).
- Hoyle R. *Pattern Formation: An Introduction to Methods*. Cambridge: Cambridge University Press (2006).
- Pecora LM, Carroll TL. Master stability functions for synchronized coupled systems. *Phys Rev Lett.* (1998) 80:4. doi: 10.1103/PhysRevLett.80.2109
- Kirk V, Silber M. A competition between heteroclinic Cycles. *Nonlinearity.* (1994) 7:1605–21. doi: 10.1088/0951-7715/7/6/005
- Krupa M, Melbourne I. Asymptotic stability of heteroclinic cycles in systems with symmetry. *Ergod Theory Dyn Syst.* (1995) 15:121–47. doi: 10.1017/S0143385700008270

AUTHOR CONTRIBUTIONS

MV and HM-O designed the study, discussed the results, and contributed to editing the manuscript. MV conducted numerical simulations, carried out the analysis, and prepared the manuscript. HM-O supervised the study.

FUNDING

Funding by the German Research Foundation (DFG, contract ME-1332/28-1) and by Jacobs University Bremen is gratefully acknowledged.

ACKNOWLEDGMENTS

We thank Ulrich Parlitz for helpful discussions during his visit at Jacobs University Bremen in summer 2019.

32. Van Kampen NG. *Stochastic Processes in Physics and Chemistry*. Amsterdam: Elsevier (1992).
33. May RM, Leonard WJ. Nonlinear aspects of competition between three species. *SIAM J Appl Math.* (1975) **29**:243–53. doi: 10.1137/0129022
34. Busse FH, Heikes KE. Convection in a rotating layer: a simple case of turbulence. *Science.* (1980) **208**:173–5. doi: 10.1126/science.208.4440.173
35. Guckenheimer J, Holmes P. Structurally stable heteroclinic cycles. *Math Proc Camb Philos Soc.* (1988) **103**:189–92. doi: 10.1017/S0305004100064732
36. Voit M, Meyer-Ortmanns H. A hierarchical heteroclinic network: controlling the time evolution along its paths. *EPJ ST.* (2018) **227**:1101–15. doi: 10.1140/epjst/e2018-800040-x

Conflict of Interest: The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2019 Voit and Meyer-Ortmanns. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.



Metrics of Emergence, Self-Organization, and Complexity for EWOM Research

Juan C. Correa*

Faculty of Psychology, Fundación Universitaria Konrad Lorenz, Bogotá, Colombia

In a recent round table organized by the Santa Fe Institute, the complexity of commerce captured the attention of those interested in understanding how complex systems science can be applicable for settings where consumers and providers interact. Despite the usefulness of applied complexity for commerce-related phenomena, few works have attempted to provide insightful ideas. This mini-review aims at providing a succinct discussion of how the metrics of emergence, self-organization, and complexity might benefit the research agenda of applied complexity and commerce/consumer studies. In particular, the paper argues possible pragmatic ways to understanding the valuable information present in word-of-mouth data found on electronic commerce platforms.

Keywords: emergence, self-organization, applied complexity, commerce-consumer research, electronic word-of-mouth

OPEN ACCESS

Edited by:

Carlos Gershenson,
National Autonomous University of
Mexico, Mexico

Reviewed by:

Diego R. Amancio,
University of São Paulo, Brazil
Oliver López-Corona,
National Council of Science and
Technology (CONACYT), Mexico

*Correspondence:

Juan C. Correa
juanc.correan@konradlorenz.edu.co

Specialty section:

This article was submitted to
Interdisciplinary Physics,
a section of the journal
Frontiers in Physics

Received: 31 October 2019

Accepted: 05 February 2020

Published: 21 February 2020

Citation:

Correa JC (2020) Metrics of
Emergence, Self-Organization, and
Complexity for EWOM Research.
Front. Phys. 8:35.
doi: 10.3389/fphy.2020.00035

1. INTRODUCTION

Emergence, self-organization, and complexity are three fundamental concepts in complex systems science [1, 2]. Nonetheless, the application of these concepts to the understanding of human behavior in the realm of commerce/consumer studies is far from being well-understood. In fact, on September 12, 2019, the Santa Fe Institute organized a discussion of this topic (https://wiki.santafe.edu/index.php/Complexity_of_Commerce_Agenda). As an extension of this matter, a well-served exposition would consist of providing a discussion on how the metrics of emergence, self-organization, and complexity [3] might benefit the research agenda of applied complexity and commerce/consumer studies. A warning note, however, should be stated beforehand. As commerce/consumer research is wide enough to be considered in one single paper, this circumstance demands the choice of a particular phenomenon. Accordingly, the remaining of this paper focuses on the consumers' "electronic word-of-mouth" (EWOM). Word-of-mouth [4] takes place when customers produce informal communications directed at other consumers about the ownership, usage, or characteristics of particular goods and services. When these communications are produced and shared through social media or electronic platforms, they are also known as "electronic word-of-mouth" [5].

Although the analysis of EWOM through statistical techniques is well-known in behavioral sciences [6, 7], the application of concepts coming from the framework of applied complexity is less frequent in the literature, being the works of Reingen and Kernan [8] and Jun et al. [9] two remarkable exceptions. Mathematical modeling or computerized simulations are also available from sociophysics [10, 11] by analyzing synthetic data.

A related yet different approach is the conceptual discussion provided here, which elaborates upon the idea of collecting natural EWOM data, preprocessing, and transform it as network data to calculate the emergence, self-organization, and complexity of its network structure. To achieve

this goal, the organization of this mini-review is as follows. The idea of EWOM as a case study for applied complexity is present in the next section which also illustrates the computational steps to follow for collecting and preprocessing EWOM data and transform it as network data. Such illustration is not an analytical coverage. Also in section 2, is present the formalization of emergence, self-organization, and complexity, by summarizing the ideas of previous works [3, 12]. Section 3, then, enumerates possible benefits and challenges for applied researchers. In section 4, the paper closes presenting possible research questions that could be used for guiding empirical studies focusing on EWOM from the perspective of applied complexity.

2. EWOM AS A CASE STUDY FOR APPLIED COMPLEXITY

From a data science perspective [13], EWOM data are not intrinsically structured, and it demands the application of natural language processing and text mining techniques [14, 15] to structure them following principles of tidy data [16]. The utility of tidying up this data lies in the possibility to leverage information mechanics (e.g., production, storage, and transmission) to gain insights into essential phenomena, such as customer engagement in online reviews [17], or quantifying the effect of online consumer reviews on new product sales [18].

In online food delivery platforms [19], it might be interesting, for example, to know possible differences among customers' experiences when consuming products of globalized fast-food chains. Because the preparation of each product follows a standardized industrial procedure in each of these globalized restaurants, several research projects can be conducted. One of these projects, for example, could be the empirical validation of agent-based models focusing on word-of-mouth dynamics with information seeking [20]. In projects of this sort, it might be revealing the description of how the dynamics of customers' positive word-of-mouth differ from the dynamics of negative word-of-mouth. With web scraping techniques for collecting real data from different globalized platforms, the possibility to characterize complaints vs. recommendations, and the estimation of customers' cultural customs when they recommend outstanding products, are certainly two other fruitful ventures. If applied researchers wish to turn their attention to customers' word-of-mouth semantics, the use of text-network analyses [21], based on principles of social network analysis [22, 23] might provide exciting answers. Working with these topics might be fruitful for those who acknowledge the imperfect nature of real-world data and yet wish to use it for theoretical development. Arguably, a brief description of how to collect and preprocess EWOM data might be illustrative for applied researchers.

2.1. Collecting and Pre-processing of EWOM Data

The use of web scraping techniques [24] is a convenient means for collecting EWOM data from online food delivery platforms. Web scraping refers to the process of extracting data from

websites automatically. The specifics on how web scraping works are beyond the scope of this paper, but the preprocessing of EWOM data deserves some mention. By its nature, EWOM data is not structured, but a convenient way to structure it is to transform customers' comments into a document-term matrix [15], whose entries show the frequency of appearance of every single word in each comment (i.e., words are arranged as rows, while comments are arranged as columns). As the number of comments generally exceeds the number of unique words that customers use for expressing their experiences, the resulting dimensionality of this matrix makes it equivalent to an incidence matrix [22]. This document-term matrix can then be re-expressed as a similarity matrix whose entries show the Jaccard index that quantifies the similarity between every word-comment unit [25]. The calculation of the Jaccard index here allows appreciating subtle semantic differences in customers' comments (e.g., a strong recommendation without hesitation on any aspect of the service vs. a recommendation accompanied by a warning regarding food variety). The knowledge of these semantic differences proves to be important for estimating the number of states for EWOM data. As these states are related to the concepts of emergence, self-organization, and complexity, it is convenient to describe them.

2.2. Emergence, Self-Organization, and Complexity of EWOM Data

In a recent paper, Santamaría-Bonfil et al. [3] summarized both the discrete and continuous measures of emergence (E), self-organization (S), and complexity (C) which are applicable to any dataset or probability distributions [12], and rely on Shannon's information theory, as pioneered by the Santa Fe Institute [1]. A few ideas about the implications of using these concepts for analyzing EWOM data are necessary at this point. The first idea posits that EWOM is a dynamic property of an open system composed of customers and sellers that interact by using an electronic platform. The second idea states the possibility of analyzing EWOM at different scales. While from a microscopic scale, one would see a series of written characters (i.e., letters, emojis, words) with a particular frequency distribution, from a macroscopic scale, one would see a set of possible semantic states (i.e., complaint, recommendation, or suggestion). As these semantic states are not trivially detectable at a microscopic level, the coordinated production of written characters allows the emergence of new behaviors (e.g. satisfied vs. unsatisfied customers, and successful vs. non-successful restaurants in the online food delivery platform). This idea is compatible with that of emergence [12] that refers to properties of a phenomenon which are present at one scale (e.g., a satisfied client) and are not at another scale (e.g., the words written by a client). According to Santamaría-Bonfil et al. [3], the concept of emergence (E) for discrete probability distribution measures the average ratio of uncertainty a process produces by new information that is a consequence of changes dynamics or scale. For continuous distributions, the interpretation of E is constrained to the average uncertainty a process produces under a specific set of statistical parameters, such as the standard deviation in a normal

distribution. The discrete and continuous versions of E are defined as

$$E_D = -K \sum_{i=1}^N p_i \log_2 p_i \quad (1)$$

$$E_C = -K(\lim_{\Delta \rightarrow 0} H(X^\Delta) + \log_2(\Delta)) \quad (2)$$

Equation (1) defines discrete E , where $p_i = P(X = x)$ is the probability of element i . Equation (2) defines continuous E , where X^Δ corresponds to discretized version of X , and Δ is the integration step, and K is a normalizing constant that constrains E in the range $[0 \leq E \leq 1]$ and is estimated as

$$K = \frac{1}{\log_2(b)} \quad (3)$$

where b corresponds to the number of bins of a probability mass function, or, in the continuous case, to the states that satisfies $P(x_i) > 0$ (i.e., recommendations, complaints, or suggestions). In addition, $\log_2(b)$ represents the maximum entropy for a distribution function with alphabet size of b (i.e., the number of characters used by customers when writing their comments). Thus, E can be deemed as the ratio between the entropy for given empirical distribution $H(X)$, and the maximum entropy for the same alphabet size $H(U)$. Now, let's turn the attention to self-organization. According to Fernández et al. [12] self-organization (S) is related to an increase in order or a reduction of entropy. Put it differently, as emergence supposes an increase of information, S should be anti-correlated with E , and this is formally expressed as

$$S = 1 - E = 1 - \left(\frac{H(P(X))}{H(U)} \right) \quad (4)$$

The numerical result of Equation (4) is also in the range $[0 \leq S \leq 1]$. With this final result, we can now realize the notion of complexity. Here, complexity represents a balance between change and regularity, allowing EWOM to adapt to contextual contingencies (e.g., showing dynamic changes as a function of the service quality of food providers or the increasing competitiveness among restaurants). While the regularity ensures the survival of information (e.g., a systematic positive opinion), change leads to the exploration of new possibilities (e.g., the emergence of recommendations for new products or services); that is, complexity describes the behavior of a system as the average uncertainty produced by emergent and regular global patterns as described by its probability distribution [3], which is formally expressed as follows:

$$C = 4 \times E \times S \quad (5)$$

In Equation (5), C is maximal when E equals S , and the highest value of C is achieved when one (or just a few) of the states is highly probable. C becomes zero when all of the states share the same probability of occurrence. The pragmatic interpretation of these metrics derive from a perspective called “*the world as evolving information*” [26]. An essential ingredient of this perspective posits the benefits of describing energy, matter, life

and cognition in terms of information. These benefits neither deny the utility of physics for describing physical phenomena, nor chemistry for chemical events, nor biology for life-related facts. Nonetheless, this perspective is meant only for the cases when the approaches of physics, chemistry, or biology are not sufficient for comprising phenomena with manifestations at different scales. The eight tentative laws of information proposed by Gershenson [26] are useful for understanding the benefits of employing the concepts of emergence, self-organization, and complexity for EWOM research.

3. ENUMERATING BENEFITS AND CHALLENGES

The recognition that EWOM is an emergent dynamic property of an open system composed of customers and sellers that interact by using an electronic platform is admittedly compatible with the idea that it changes as time goes by; i.e., *the law of information transformation* as proposed in *the world as evolving information*. As any customer can perceive (i.e., read) the information provided by other customers regarding their experiences in dealing with a particular seller, this sort of customer-to-customer interaction is also compatible with *the law of information propagation*. If, for example, the seller-to-customer interaction preserves itself in terms of a systematic presence of customers' complaints (i.e., one of the probable semantic states), this circumstance opens the possibility for the electronic platform to penalizing the seller (e.g., when Amazon automatically returns the money paid by the customer after reporting any irregularity with the quality of the product shipped by the seller). In this last case, the so-called *law of requisite complexity* would be taking place, resulting from the platform and the seller. The ability of a seller to generate the best service possible so as to create a critical balance between a stable positive EWOM with a rather minimum amount of negative EWOM, would be deemed as *the law of information criticality*. If we accept the idea that EWOM is a powerful online information source that influences online shopping [27], then we can realize that this information is having a certain control over its environment, which conforms to *the law of information organization*. *The law of information self-organization* stating that information tends to its preferred, most probable state also has an implication for EWOM studies. Because customers engage in the so-called “collaborative consumption” [19], the publication of opinions aiming at persuading other's decisions will create the possibility of a shared and dominant opinion regarding seller's conduct. This fact also relates to *the law of information potentiality*, according to which a customer can give different potential meanings to information. Finally, *the law of information perception* implies that the perception of customers might be generalized so as to respond to novel information. Even though the precise situation and context are always unique, this creates some sort of uncertainty, and this is intrinsically related to Shannon's entropy, as explained by Fernández et al. [12]. This last concept permits me to enumerate some challenges for

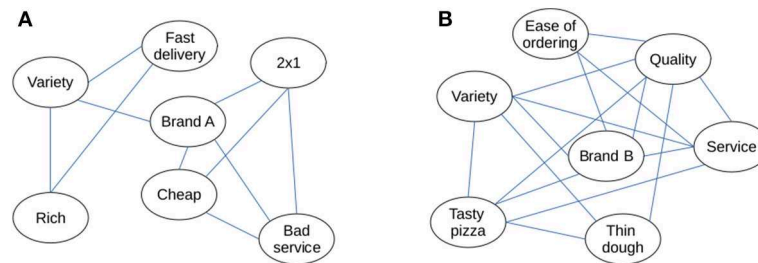


FIGURE 1 | Two consumer associative networks resulting from text-network analysis. **(A)** Shows a hypothesized network of customers' comments for brand A. **(B)** Shows a hypothesized network of customers' comments for brand B.

researchers who acknowledge the narrowness in the scope, lenses, and epistemology of their discipline [28].

The first challenge for researchers with little or no knowledge of applied complexity is the idea that all economic agents (i.e., consumers, sellers, and platforms) can be seen as interacting components of a dynamic system. As agents can be seen as systems too, a second challenge is the use of social network analysis [22] to understand the relationships of economic agents from a systemic viewpoint. Few works have followed this orientation without using the concepts of emergence, self-organization, and complexity [29]. For example, Henderson et al. [30] showed several empirical examples of consumer associative networks to mapping an extensive array of branding effects, including branded features, driver brands, complements, co-branding, cannibalization, brand parity, brand dilution, brand confusion, counter-brands, and segmentation. The idea of consumer associative networks proved to be essential for the so-called “goal systems theory” proposed by Kopetz et al. [31]. This theory posits that the study of the goal-action interaction, taking place in a cognitive and motivational processes of the consumer, might be revealing for understanding a set of consumer-related phenomena including product variety search, impulsive buying, preferences, choices, and regret. Rocha and Holme [32] showed another applied perspective when they studied the network organization of consumer complaints. Although the orientation of these works might be the standard for scientific associations, such as the complex systems society, my own impression is that they remain widely ignored by members of other applied-oriented associations, such as the society for consumer psychology, or the association for consumer research.

The ideas mentioned above call for the development of interdisciplinary perspectives that demand the search for novel insights. For example, the concept of “antifragility” [33] might be fruitful to explain why some products become best-sellers even after receiving a bunch of negative reviews. The search for novel insights also demands the use of other tools for collecting and analyzing EWOM data. It is beyond the scope of this mini-review to provide a thorough description of these tools, but they include the use of agent-based modeling and simulation [20], web scraping, natural language processing, text mining and network analysis [34]. As these techniques are easily implemented in object-oriented programming languages,

applied researchers might regard strategic their learning. After all, these programming languages offer other benefits, such as reproducibility; allowing others to follow the computational procedures that allow them to get the same results reported in a publication [35], or scalability; employing technologies capable of collecting and analyzing massive amounts of data [36]. The goals of scientific projects, such as FutureICT [37] that promote the use of the power of information to explore social and economic life, certainly call for multidisciplinary collaboration. All that is needed is the proposal of empirical studies where commerce/consumer studies and applied complexity can meet. EWOM research from an applied complexity perspective might be deemed as one of the several cases aligned with these goals.

4. CONCLUDING REMARKS

Until this point, it should be clear how applied complexity can provide several contributions to the study of EWOM research. While the concept of consumer associative network is useful for understanding EWOM data from a psychological viewpoint [30, 31], the concepts of emergence, self-organization, and complexity have not been integrated. This integration might be better understood with an example. **Figure 1** shows two consumer associative networks resulting from the procedures described in section 2.1. Although both of these figures reveal the network structure of EWOM data for two different brands of pizzas, the network on the left shows a different structure of the network on the right.

With the calculus of *E*, *S*, and *C*, commerce/consumer researchers end up with a set of proxies to the degree of customers' comments diversification, customers' comments polarization, and the diversification-polarization balance, respectively. Because the network structure of EWOM data might change as time goes by, then a dynamic analysis of these changes might help commerce/consumer researchers understand the (external) factors that act upon these structures (e.g., How effective are promotions to increase and maintain the number of positive comments?). The comparison between these structures is another issue to explore (e.g., How similar are the network structures of EWOM data for two restaurants of a globalized fast-food chain operating in different countries?). Finally, the power of emergence, self-organization, and complexity for

predicting future sales could be another related topic (How sensitive are sales to significant changes in the network structure of EWOM data?). These topics are relevant when we consider the case of Uber Eats, Just-Eat, Food Panda, or Delivery Hero, as business models that facilitate the interaction between customers and restaurants [19]. Working with EWOM data collected from these globalized platforms turns out to be an empirical field with unknown opportunities for complex systems scientists. The reason behind this statement is the gap between theory and observation. In network analysis, for example, idealized-mathematical illustrations make use of networks with few nodes and edges, but what would happen if we need to work with vast data sets of comments for a numerous collection of food

providers? How much scalability would be required for analyzing a disproportionate set of data? These questions posit important challenges for developers of cloud computing technologies, such as Data bricks or Google Cloud.

AUTHOR CONTRIBUTIONS

JC conceived and wrote the paper.

FUNDING

This research was funded by Fundación Universitaria Konrad Lorenz under research grant number 9IN11191.

REFERENCES

- Prokopenko M, Boschetti F, Ryan AJ. An information-theoretic primer on complexity, self-organization, and emergence. *Complexity*. (2009) **15**:11–28. doi: 10.1002/cplx.20249
- Polani D, Prokopenko M, Yeager LS. Information and self-organization of behavior. *Adv Complex Syst*. (2013) **16**:1303001. doi: 10.1142/S021952591303001X
- Santamaría-Bonfil G, Gershenson C, Fernández N. A package for measuring emergence, self-organization, and complexity based on Shannon entropy. *Front Robot AI*. (2017) **4**:1–12. doi: 10.3389/frobt.2017.00010
- Berger J. Word of mouth and interpersonal communication: a review and directions for future research. *J Consum Psychol*. (2014) **24**:586–607. doi: 10.1016/j.jcps.2014.05.002
- Chen M, Chen J, Xue W. Research on the influence mechanism of eWOM on selection of tourist destinations—the intermediary role of psychological contract. In: Xu J, Ahmed SE, Cooke FL, Duca G, editors. *Proceedings of the Thirteenth International Conference on Management Science and Engineering Management*. Cham: Springer International Publishing (2019). p. 654–67.
- Godes D, Mayzlin D. Using online conversations to study word-of-mouth communication. *Market Sci*. (2004) **23**:545–60. doi: 10.1287/mksc.1040.0071
- Chevalier JA, Mayzlin D. The effect of word of mouth on sales: online book reviews. *J Market Res*. (2006) **43**:345–54. doi: 10.1509/jmkr.43.3.345
- Reingen PH, Kernan JB. Analysis of referral networks in marketing: methods and illustration. *J Market Res*. (1986) **23**:370–8. doi: 10.1177/002224378602300407
- Jun T, Kim JY, Kim BJ, Choi MY. Consumer referral in a small world network. *Soc Netw*. (2006) **28**:232–46. doi: 10.1016/j.socnet.2005.07.001
- Chakrabarti AS, Sinha S. “Hits” emerge through self-organized coordination in collective response of free agents. *Phys Rev E*. (2016) **94**:042302. doi: 10.1103/PhysRevE.94.042302
- Ishii A, Kawahata Y. Sociophysics analysis of the dynamics of peoples’ interests in society. *Front Phys*. (2018) **6**:89. doi: 10.3389/fphy.2018.00089
- Fernández N, Maldonado C, Gershenson C. Information measures of complexity, emergence, self-organization, homeostasis, and autopoiesis. In: Prokopenko M, editor. *Guided Self-Organization: Inception*. Vol. 9. Berlin; Heidelberg: Springer (2014). p. 19–51.
- Wickham H, Grolemund G. *R for Data Science: Import, Tidy, Transform, Visualize, and Model Data*. Boston, MA: O’Reilly Media, Inc. (2017).
- Silge J, Robinson D. Tidytext: text mining and analysis using tidy data principles in R. *J Open Source Softw*. (2016) **1**:37. doi: 10.21105/joss.00037
- Benoit K, Watanabe K, Wang H, Nulty P, Obeng A, Müller S, et al. Quanteda: an R package for the quantitative analysis of textual data. *J Open Source Softw*. (2018) **3**:774. doi: 10.21105/joss.00774
- Wickham H. Tidy data. *J Stat Softw*. (2014) **59**:1–23. doi: 10.18637/jss.v059.i10
- Thakur R. Customer engagement and online reviews. *J Retail Consum Serv*. (2018) **41**:48–59. doi: 10.1016/j.jretconser.2017.11.002
- Cui G, Lui HK, Guo X. The effect of online consumer reviews on new product sales. *Int J Electron Commerce*. (2012) **17**:39–58. doi: 10.2753/JEC1086-4415170102
- Correa JC, Garzón W, Brooker P, Sakarkar G, Carranza SA, Yunado L, et al. Evaluation of collaborative consumption of food delivery services through web mining techniques. *J Retail Consum Serv*. (2019) **46**:45–50. doi: 10.1016/j.jretconser.2018.05.002
- Thiriot S. Word-of-mouth dynamics with information seeking: information is not (only) epidemics. *Phys A*. (2018) **492**:418–30. doi: 10.1016/j.physa.2017.09.056
- Choi Y, Kweon SH. A semantic network analysis of the newspaper articles on big data. *J Cybercommun Acad Soc*. (2014) **31**:241–86.
- Wasserman S, Faust K. *Social Network Analysis: Methods and Applications*. Cambridge: Cambridge University Press (1994).
- Barabasi AL. *Linked: The New Science of Networks*. New York, NY: Perseus Books Group (2002).
- Munzert S, Rubba C, Meißner P, Nyhuis D. *Automated Data Collection With R: A Practical Guide to Web Scraping and Text Mining*. New York, NY: John Wiley & Sons (2014).
- Sun J, Tang J. A survey of models and algorithms for social influence analysis. In: Aggarwal CC, editor. *Social Network Data Analytics*. New York, NY: Springer (2011). p. 177–214.
- Gershenson C. The world as evolving information. In: Minai A, Braha D, Bar-Yam Y, editors. *Unifying Themes in Complex Systems*. Vol. VII. Heidelberg: Springer (2012). p. 100–15.
- Bag S, Tiwari MK, Chan FT. Predicting the consumer’s purchase intention of durable goods: an attribute-level analysis. *J Bus Res*. (2019) **94**:408–19. doi: 10.1016/j.jbusres.2017.11.031
- Pham MT. The seven sins of consumer psychology. *J Consum Psychol*. (2013) **23**:411–23. doi: 10.1016/j.jcps.2013.07.004
- Seoane LF, Solé R. The morphospace of language networks. *Sci Rep*. (2018) **8**:10465. doi: 10.1038/s41598-018-28820-0
- Henderson GR, Iacobucci D, Calder BJ. Brand diagnostics: mapping branding effects using consumer associative networks. *Eur J Oper Res*. (1998) **111**:306–27. doi: 10.1016/S0377-2217(98)00151-9
- Kopetz CE, Kruglanski AW, Arens ZG, Etkin J, Johnson HM. The dynamics of consumer behavior: a goal systemic perspective. *J Consum Psychol*. (2012) **22**:208–23. doi: 10.1016/j.jcps.2011.03.001
- Rocha LEC, Holme P. The network organisation of consumer complaints. *Europhys Lett*. (2010) **91**:28005. doi: 10.1209/0295-5075/91/28005
- Pineda OK, Kim H, Gershenson C. A novel antifragility measure based on satisfaction and its application to random and biological Boolean networks. *Complexity*. (2019) **2019**:3728621. doi: 10.1155/2019/3728621
- Amancio DR, Silva FN, Costa LdF. Concentric network symmetry grasps authors’ styles in word adjacency networks. *Europhys Lett*. (2015) **110**:68001. doi: 10.1209/0295-5075/110/68001

35. McNutt M. Reproducibility. *Science*. (2014) **343**:229. doi: 10.1126/science.1250475
36. Bakshi S, Jagadev AK, Dehuri S, Wang GN. Enhancing scalability and accuracy of recommendation systems using unsupervised learning and particle swarm optimization. *Appl Soft Comput J*. (2014) **15**:21–9. doi: 10.1016/j.asoc.2013.10.018
37. San Miguel M, Johnson JH, Kertesz J, Kaski K, Díaz-Guilera A, MacKay RS, et al. Challenges in complex systems science. *Eur Phys J Spec Top*. (2012) **214**:245–71. doi: 10.1140/epjst/e2012-01694-y

Conflict of Interest: The author declares that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2020 Correa. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.



How Computation Is Helping Unravel the Dynamics of Morphogenesis

David Pastor-Escuredo^{1,2,3*} and Juan C. del Álamo^{2,3,4,5,6*}

¹ LifeD Lab, Madrid, Spain, ² Mechanical and Aerospace Engineering Department, University of California San Diego, San Diego, CA, United States, ³ Institute for Engineering in Medicine, University of California San Diego, San Diego, CA, United States, ⁴ Mechanical Engineering Department, University of Washington, Seattle, WA, United States, ⁵ Center for Cardiovascular Biology, University of Washington School of Medicine, Seattle, WA, United States, ⁶ Institute for Stem Cell and Regenerative Medicine, University of Washington School of Medicine, Seattle, WA, United States

OPEN ACCESS

Edited by:

Carlos Gershenson,
National Autonomous University of
Mexico, Mexico

Reviewed by:

Ignazio Licata,
Institute for Scientific Methodology
(ISEM), Italy
Reinaldo Roberto Rosa,
National Institute of Space Research
(INPE), Brazil

*Correspondence:

David Pastor-Escuredo
david@lifedlab.org
Juan C. del Álamo
jc@ucsd.edu;
juancar@uw.edu

Specialty section:

This article was submitted to
Interdisciplinary Physics,
a section of the journal
Frontiers in Physics

Received: 31 October 2019

Accepted: 04 February 2020

Published: 28 February 2020

Citation:

Pastor-Escuredo D and del Álamo JC
(2020) How Computation Is Helping
Unravel the Dynamics of
Morphogenesis. *Front. Phys.* 8:31.
doi: 10.3389/fphy.2020.00031

The growing availability of imaging data, calculation power, and algorithm sophistication are transforming the study of morphogenesis into a computation-driven discipline. In parallel, it is accepted that mechanics plays a role in many of the processes determining the cell fate map, providing further opportunities for modeling and simulation. We provide a perspective of this integrative field, discussing recent advances and outstanding challenges to understand the determination of the fate map. At the basis, high-resolution microscopy and image processing provide digital representations of embryos that facilitate quantifying their mechanics with computational methods. Moreover, innovations in *in-vivo* sensing and tissue manipulation can now characterize cell-scale processes to feed larger-scale representations. A variety of mechanical formalisms have been proposed to model cellular biophysics and its links with biochemical and genetic factors. However, there are still limitations derived from the dynamic nature of embryonic tissue and its spatio-temporal heterogeneity. Also, the increasing complexity and variety of implementations make it difficult to harmonize and cross-validate models. The solution to these challenges will likely require integrating novel *in vivo* measurements of embryonic biomechanics into the models. Machine Learning has great potential to classify spatio-temporally connected groups of cells with similar dynamics. Emerging Deep Learning architectures facilitate the discovery of causal links and are becoming transparent and interpretable. We anticipate these new tools will lead to multi-scale models with the necessary accuracy and flexibility to formulate hypotheses for *in-vivo* and *in-silico* testing. These methods have promising applications for tissue engineering, identification of therapeutic targets, and synthetic life.

Keywords: morphogenesis, cell mechanics, multi-scale modeling, morphomechanical fields, deep learning, cell fate map, fluorescence microscopy, digital embryo

INTRODUCTION

Embryogenesis is a complicated ensemble of processes by which a single cell turns into a multi-cellular living organism. Through various developmental stages, the cell population proliferates while tissues develop, change their properties, differentiate, and gain their specific functionality [1]. During embryogenesis, biochemical, genetic, and epigenetic factors interact, forming a tangled network of processes with diverse physical length scales and time scales [2, 3]. Remarkably,

the robustness and variability of these processes are balanced to make possible the reproducibility and diversity of living specimens [4].

Mechanics plays a central role in shaping the embryo [5, 6]. Gene expression gradients regulate tissue patterning and cellular properties, such as rheology, adhesion, and contractility [7]. At the same time, the embryonic cells sense mechanical cues from their microenvironment and convert them to biochemical signals, including gene expression [8, 9]. These cues are essential to guide morphogenesis but also tissue repair, given that immature cells can retain significant plasticity and reprogram in response to external forces [10, 11]. The cross-talk between biophysical and biochemical processes involves multiple mechanisms and molecules and occurs in multiple scales [12]. Besides, cells can follow complex trajectories within the developing embryo, thereby creating and being exposed to continuous changes in the microenvironment [13–16].

Researchers have been long interested in discovering mechanistic links between physical processes and gene expression that lead to cell fate determination [17–21]. Recent advances in microscopy, modeling, and computation have enabled quantifying 2D and 3D mechanical forces and rheological properties in multi-cellular colonies, including live developing embryos [22–28]. These methods provide local data in space and time, and analyzing them to unravel cell fate maps is challenging. High-resolution, long-term observation in two or three dimensions is desirable to consider the whole range of scales at which mechanics can impact cell fate. Still, it complicates the analysis further because it involves massive amounts of data. Furthermore, the statistical treatment of the data needs to accommodate the highly heterogeneous and time-evolving properties of developing tissues [29–31].

This perspective discusses current advances in computational methods for the characterization of mechanical processes during embryogenesis and how these processes influence cell fate. Sections Digital Reconstruction of Embryogenesis, *In vivo* Quantification of Forces and Mechanical Properties, and Computational Models in this perspective are organized according to key steps in the analysis of experimental data and relevant methodological approaches. Each section presents our view into key advances and outstanding challenges. Section Morphomechanical Domains in Developing Tissues: Follow the Cell, Not the Voxel proposes a paradigm to deal with the massive data produced by experimental techniques and construct a multi-scale representation of embryo dynamics. Finally, section Understanding Multi-Scale Embryonic Dynamics by Machine Learning presents problems at the intersection between morphogenesis and Machine Learning that has not been so far tackled by the community.

DIGITAL RECONSTRUCTION OF EMBRYOGENESIS

Progress in live microscopy and fluorescence reporters now allow high-resolution, time-lapse imaging of developing

embryos in two and three dimensions [32–35]. Image analysis and computer vision methods can now create digital atlases of developing embryos (**Figures 1A–F**). These atlases contain spatio-temporal information about cell and tissue morphology, cell lineages, and functional patterns, such as gene expression or protein activity [36–42]. Moreover, novel visualization tools allow for systematically browsing these digital embryos (**Figure 1**), and integrating them into numerical simulations and machine learning algorithms [36, 42, 43].

The three-dimensional *in-vivo* imaging of whole embryos has challenges associated with image resolution, quality, and artifacts (e.g., anisotropic point spread function). Besides, photobleaching and phototoxicity make it challenging to extend imaging over intervals long enough to capture relevant morphogenetic processes. Multi-view light-sheet microscopy (LSM) [35, 44, 45] and view fusion algorithms [46, 47] allow for 3D imaging large embryos with cellular isotropic resolution. Recently, advances proposing adaptive optics and lattice LSM with ultrathin light-sheet excitation featured, promising sub-cellular resolution during long-term observation [48].

Reconstructing the multi-scale dynamics of embryogenesis requires not only long-term imaging with sub-cellular spatial resolution but also sub-minute temporal resolution. An established approach to achieve these joint demands is to record images of several embryos within the same cohort with different temporal resolutions and to register the resulting images spatially onto a common template [49, 50]. The projected growth in computing power of microscopy systems (e.g., by embedded GPU computing) makes it possible to envision enhanced microscopes with real-time adaptive multi-scale observation [51, 52].

Image processing workflows must be able to handle the massive amounts of complex data resulting from microscopy modalities to provide a quantification of structures, motion, and hierarchy [3, 53]. Intensity-based methods, such as optical flow or image registration produce continuous velocity fields [53–55] that can leverage the powerful modeling and descriptive tools of continuum and statistical mechanics [56, 57]. On the other hand, tracking the motion and divisions of single cells yields discrete cell lineages, which presents apparent advantages [4, 36, 40].

Motion estimation is critical because determining cell fate involves reconstructing 3D cell trajectories across the various developmental stages, imposing quasi-error-free requirements (**Figures 1E,F**). Deep Learning tools, such as Convolutional Neural Networks can help to improve the performance under challenging conditions, such as deep-tissue segmentation provided tagged training data [58, 59]. Interactive annotation tools for correction and validation are still a suitable approach to generate reliable expert-driven data [36, 42, 43, 60] and potentially allow crowdsourced results [61]. Beyond image data repositories, sharing detailed experiment metadata through systematic frameworks (e.g., based on ontologies) can provide a “Big Data” substrate for machine learning to optimize pipelines.

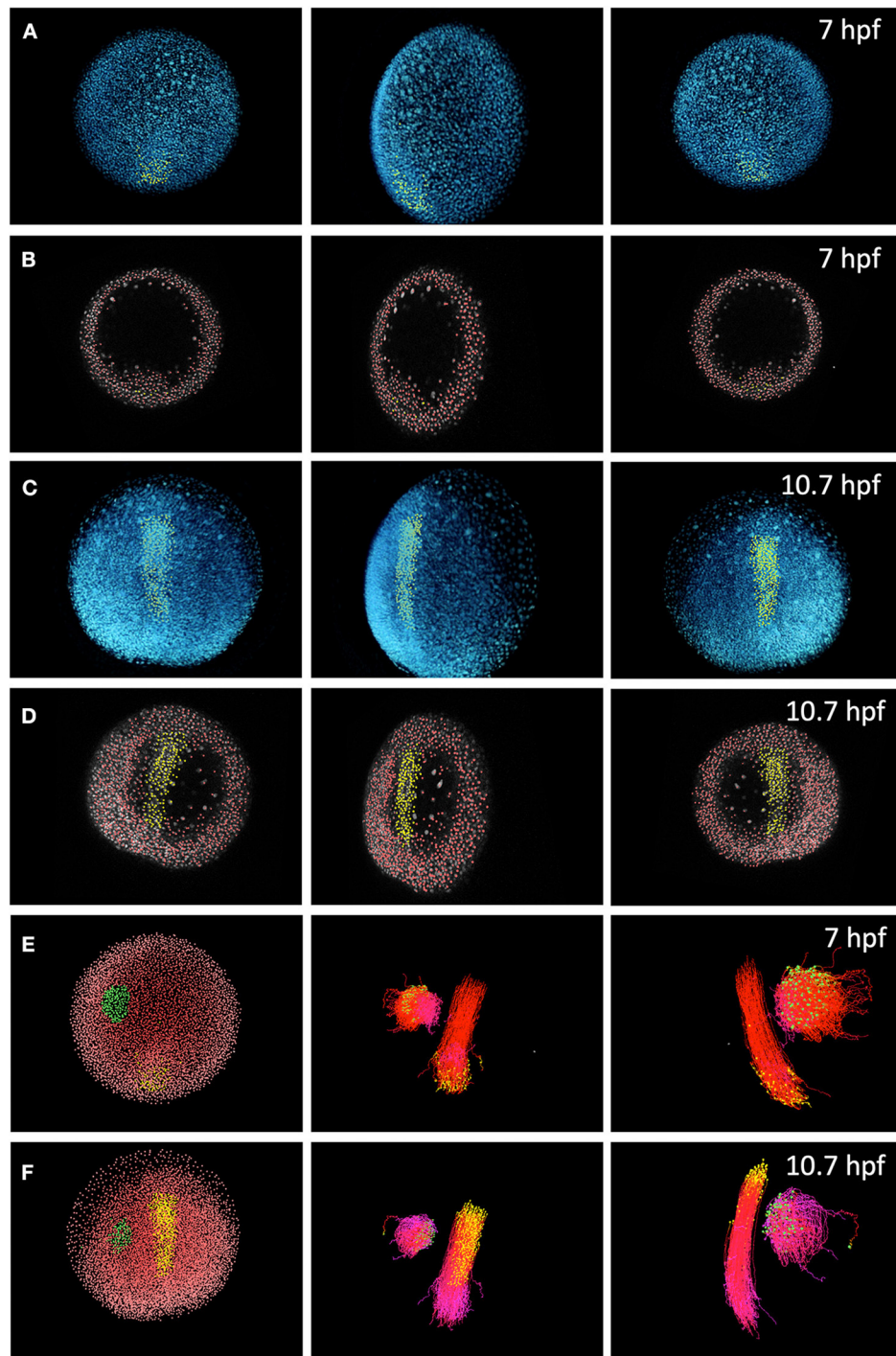


FIGURE 1 | Digital reconstruction of cell trajectories in a zebrafish embryo. **(A)** Three-dimensional (3D) rendering of cell nuclei (blue) in raw data and selection of cells (yellow) in the dorsal line, performed at 10.7 h-post-fertilization and backtracked to 7 hpf. From left to right, each panel shows a different spatial orientation (animal pole, lateral and ventral). **(B)** Detected cell nuclei (red) and cell selection as in **(A)** (yellow), shown in a spatial slice over the orthoslice of the raw data channel (gray). Same time step and view angles as in **(A)**. **(C)** 3D rendering of cell nuclei (blue) and selection of cells of the dorsal line at 10.7 hpf (yellow). **(D)** Cell detection (red) and cell selection (yellow) in the same slice as **(B)**. **(E)** *Left*: Two cell selections (green and yellow) over nuclei detection domains at 7 hpf. *Middle*: forward tracking (yellow to red colormap indicates time advancement) of the two selected cell domains. *Right*: forward tracking of the two cell selections from a lateral point of view. **(F)** Cell selections (green and yellow) at 10.7 hpf. Cell backward tracking (yellow to purple colormap) in same view angles than **(E)**. This dataset and the visualization tool Mov-IT are freely available from the BioEmergences open workflow <http://bioemergences.iscpif.fr/bioemergences/openworkflow-index.php> [36].

IN VIVO QUANTIFICATION OF FORCES AND MECHANICAL PROPERTIES

Digital reconstruction of morphogenesis already provides structured data, but embedding biophysical formalisms is invaluable to decipher multi-scale dynamics. The study of biophysics in single cells is not new: the measurement of the rheological properties of cells, their internal stresses and the forces they exert has received significant attention in the past two decades [12, 23, 62]. The requirement of non-invasiveness, three-dimensionality, and the need for calibrated sensors that sensitive enough to resolve minute forces and deformations make this task particularly challenging in live embryos.

Laser ablation was one of the pioneering methods to quantify embryonic mechanics *in vivo*. This technique produces a localized cut in a tissue, which allows for estimating tissue tensions by letting the ablated region relax to a stress-free configuration [63–65]. This technique is still widely used but it is disruptive. A non-invasive alternative is to use fluorescence reporters to measure acto-myosin activity as a surrogate metric of force generation. Still, both methods rely on independent measurements of tissue rheology [66, 67]. Molecular sensors based on fluorescence resonance energy transfer (FRET) also provide a minimally invasive means of measuring forces *in vivo* [68]. This modality is very attractive since it probes the tension born by specific molecules. However, it requires careful calibration, does not provide vector or tensor data, and needs a different sensor to measure the tension born by each molecule. It is undoubtable that these approaches will continue to shed light on numerous embryonic processes. Even so, their critical examination has kindled the search for easy-to-calibrate quantifications of the strains, stresses and material properties inside live tissues.

Because *in vitro* assays allow for careful control of experimental parameters, they have experienced significant progress in the past 20 years, thus offering valuable lessons for the development of *in vivo* techniques. In particular, there is a diversity of force microscopy methods that exploit the linear properties and high deformability of hydrogels to provide sensitive, calibrated strain-stress sensors. Cells are cultured on these hydrogels, the deformation caused by the cells on the hydrogel is measured, typically by tracking the motion of fiduciary markers (e.g., fluorescent microspheres), and the traction forces exerted by the cells are recovered from the measured deformations [69–71]. Monolayer Stress Microscopy is an extension of traction microscopy that quantifies the collective distribution of intracellular stress in thin confluent cell cultures [72]. A similar approach was proposed to estimate ventral furrow invagination in *Drosophila* although in that case the stress-free configuration was not known [73]. Of note, traction forces can be highly three-dimensional even when the cells are plated on flat hydrogels [74], leading to significant bending and additional intracellular stress in cell monolayers [75]. Quantifying the

forces involved in epithelial bending and invagination could offer new biomechanical insights about the morphogenesis of tissues and organs.

In live developing embryos, it is now feasible to measure strains (and strain rates, **Figure 2**) at the cellular level by tracking the morphological changes of segmented cells [55, 76]. Tissue-level strain fields can be derived from cell tracking and optical flow methods (**Figures 2A–C**) [57, 77]. By combining the cell-level and tissue-level strain quantifications it is possible to infer tissue rearrangements, such as cell deformation and cell intercalation [55, 60, 77, 78]. These metrics can be overlaid with functional data, such as gene expression and acto-myosin activity, to provide a correlation-based understanding of tissue dynamics [53, 77, 79]. Moreover, continuum strain fields enable the quantification of internal stresses based on a prescribed mechanical model for the embryo. These formulations are very advantageous—they allow for writing sets of equations that can be solved analytically or numerically to fully characterize the mechanical state of the tissue [80]. A mechanical formalism that has been applied to developing embryos with notable success relies on enforcing static equilibrium of forces between intracellular pressure and cortical tension. This formulation leads to a geometrical problem for cell shapes that can be closed by analyzing experimental images [81–86]. However, it must be recalled that embryonic tissue properties are heterogeneous, highly non-linear and time-evolving, which makes it challenging to develop mechanical formalisms that are uniformly valid across different regions of space, instants of time, and genetic and pharmacological manipulations. Furthermore, a significant challenge is to establish the stress-free reference state to properly quantify visco-elastic forces.

A recent approach for the *in-vivo* characterization of embryonic mechanics, without prior assumptions, consists of injecting microdroplets or hydrogel microspheres of size comparable to one cell, and that can act as calibrated sensors and/or actuators (**Figure 2**) [87]. After appropriate functionalization by surface coating, these sensors can be made biocompatible and are internalized by the embryo, thereby minimizing the invasiveness of the method. Incompressible fluorescent oil-droplets allow for quantifying anisotropic stresses [88], whereas hydrogel droplets with characterized compressibility allow for quantifying isotropic ones [89]. Moreover, ferrofluid droplets can be act as active sensors to measure the local tissue rheology [90]. An additional feature of these sensors is that they move with their neighboring cells during development, thus providing valuable information about the temporal evolution of mechanical stresses and tissue rheology. Their limitations stem from reduced sampling ability, given by the limited number of sensors that can be used per embryo, and the current lack of scalable computational frameworks to relate the measurements with cell fate determination. Even so, it is reasonable to expect that emerging innovations will simplify the implementation of these techniques, enabling their widespread application.

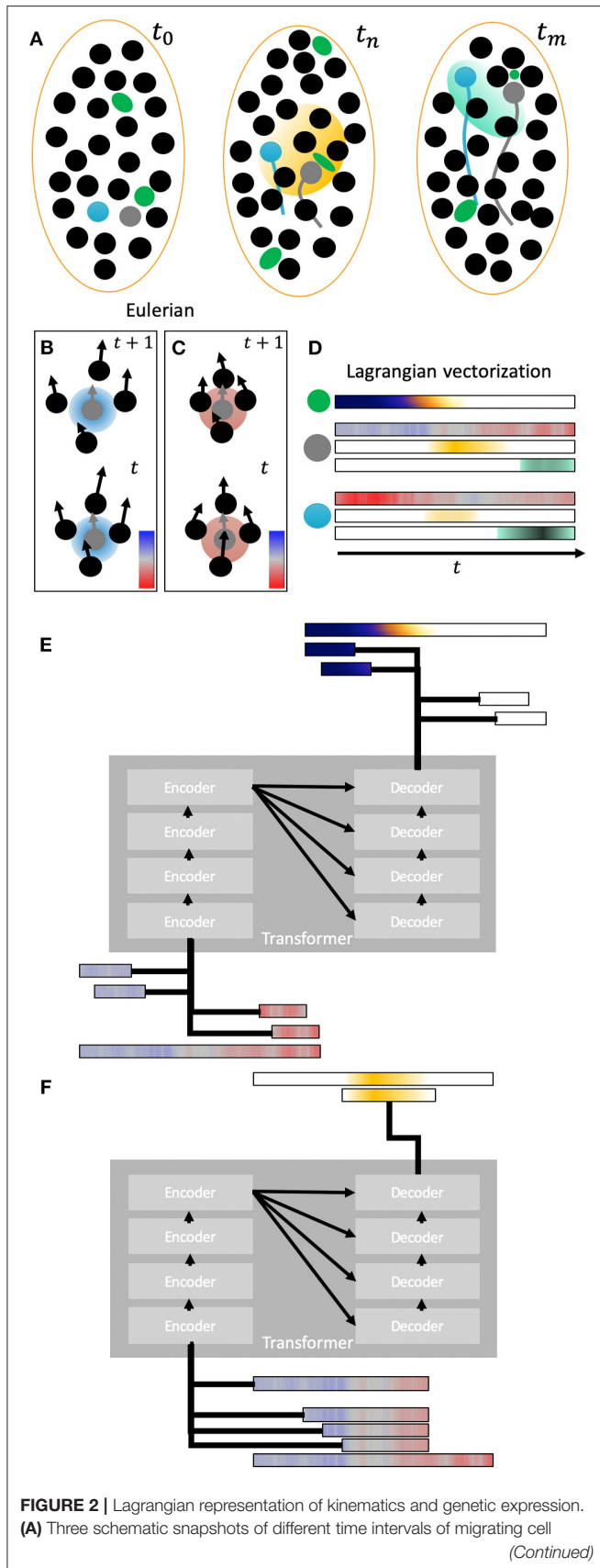


FIGURE 2 | nuclei within an area. The trajectories of two cells (gray and blue) are highlighted. In the second time step a gene expression pattern in yellow is shown affecting several cells. In the last time step another different gene expression pattern is represented in cyan affecting another set of cells. Green blobs represent mechanical sensors that sense local deformation. **(B)** Two snapshots showing the relative displacements of neighbor cells with respect of a reference cell (gray). These relative displacements are translated into a kinematic descriptor of relative area change rate that representation expansion (blue) and compression (red). The relative displacements in this schematic panel imply a local expansion (blue value) around the reference cell between timesteps t and $t + 1$ as shown by the average increasing distance between the cell nuclei. **(C)** same schematic than in **(B)** for a compressive case with cells getting closer to the reference cell (red value). **(D)** Lagrangian vectorization of compression/expansion descriptor [same colormap than **(B,C)**] and gene expression along time for the two reference cells. The data sensed with the mechanical probe is also vectorized in a Lagrangian representation with colormap dark blue to white. Gene expression is vectorized along the reference cell trajectories. **(E)** Schematic of a transformer (encoder-decoder) architecture trained to infer local forces from deformation measurements (input). The input is segmented into different temporal frames as subvectors. **(F)** Schema of a transformer architecture trained to infer mechanical factors (input) involved in the appearance of expression patterns at different temporal scales.

COMPUTATIONAL MODELS

Computational models with explanatory and predictive power can infer causal links and contribute to the mechanistic understanding of embryogenesis. These models allow researchers to observe processes, reverse engineer mechanisms, and test hypotheses with much looser constraints than pharmacological or genetic manipulations. Many biological problems involving collective cell-cell and cell-matrix interactions have been simulated using discrete, continuum, and hybrid physical models [91, 92]. Discrete agent-based models initially considered cellular movements within a lattice to investigate multicellular interactions [93]. Lattice-free agent-based models consider continuous movements of each agent. A common approach is to conceptualize cells as objects with fixed geometry and biophysical properties, whose trajectories are dictated by the balance of forces exerted by their neighbors and the environment [94]. Subcellular resolution can be achieved through agent-based models in which each agent is deformable and occupies several nodes [95]. The cellular Potts model (CPM) is an energy-based stochastic model, typically defined on a lattice that can have subcellular resolution, that is particularly well-suited to deal with large deformations and multi-scale phenomena [96]. These features make the CPM well-suited to simulate collective cell dynamics in a diversity of scenarios, including morphogenesis [97]. While they are mostly phenomenological, these models are a promising, computationally efficient approach to study how meso-scale multicellular phenomena emerge from the self-organization of sub-cellular and cellular processes.

The cellular Potts model was initially applied to quantify epithelial dynamics including the rearrangements of different cells [98]. Subsequently, the CPM has provided insight about how cortical tension and cell adhesions drive cell sorting and tissue organization [99, 100]. More recently, agent-based

models have proven useful to integrate mechanical cues with gene expression. Epithelial and mesenchymal tectonics were simulated together with gene regulatory network dynamics to recapitulate the dynamics of early zebrafish development [26]. Deformable agent-based models are a promising approach to quantify mechanotransduction, the heterogeneity of embryonic tissues, and their impact in larger-scale developmental processes [92, 101].

Vertex models bridge the discrete and continuum descriptions [102]. In these models, each cell is approximated by a polygon in 2D or a polyhedron in 3D, and the tissue measurements are sampled at the junction of three or more cells [102, 103]. Vertex models provide more information on cell interfaces than agent-based models permitting the analysis of topological changes in the cell environment [104, 105]. Curved cell geometries can be resolved with finite-elements [106, 107], and the biophysical interaction between the membrane and the cytoplasmic fluid can also be incorporated using immersed boundary methods [108]. Vertex models have been widely applied to study the mechanics of epithelia, which are represented as manifolds that can fold or invaginate [109–114]. These models have made contributions to our understanding various tissue behaviors: growth [115–117], cell division and packing [118], planar polarity [119] and the formation of compartments [120]. Dynamic cellular finite-element models have been also proposed for individual and collective cell movements and mechanics [121].

As stated above, continuum models can adapt mechanical theories, such as hydrodynamics and statistical mechanics to live matter [122], taking advantage of a massive body of knowledge and powerful tools from applied mathematics and computation, such as stability theory, perturbation methods, and computational fluid dynamics. In addition to providing a means to relate measurements of strain fields to internal stresses [78, 123], these models are well-suited to perform predictive simulations large-scale embryo dynamics. The widely studied formation of the ventral furrow in *Drosophila* [124] is a good example of a process governed by hydrodynamics [56]. Most continuum models are limited by their inherent coarse-grained, but fusion between these models and agent-based models could help resolve the contribution individual cells to tissue behavior [125].

MORPHOMECHANICAL DOMAINS IN DEVELOPING TISSUES: FOLLOW THE CELL, NOT THE VOXEL

Although microscopy experiments provide increasingly rich data about embryonic development, the data is obtained in a form that makes it difficult to extract the relationships between cellular and subcellular dynamics, large-scale biomechanical phenomena, and cell fate maps. The root for this difficulty can be illustrated using the analogy between the cell trajectories and a flow; observation through the microscopy imposes a perspective in a fixed reference frame as an external observer of embryogenesis (i.e., Eulerian frame). However, a perspective as an internal observer that records data along the trajectory

of each cell would be more suitable (i.e., Lagrangian frame). The Lagrangian framework allows for computing deformation rates and finite deformations over arbitrarily long time intervals [57]. It also helps discover Lagrangian coherent structures [126, 127] formed by cells that experience similar histories of mechanical cues, and which potentially organize the embryogenic flow (Figure 2D).

The Lagrangian trajectories of embryonic cells can be obtained by single-cell tracking or by approximating their motion as a continuous flow [14, 15, 53]. Moreover, in the Lagrangian framework, descriptors related to morphology, mechanics, genetics, etc. can be expressed in terms of the cell trajectories at specific time intervals. The usefulness of this approach depends on whether it can identify true *morphomechanical fields*. That is, if it finds connected domains of cells with a similar history of cues, if these domains are reproducible across several specimens, and if they can be related to the fate maps. We previously showed that machine learning does identify *morphomechanical fields* by classifying cell populations with similar Lagrangian cues either via clustering or with training data [57]. Comparison of cohorts can be either performed using a canonical embryo as reference or computing a statistical average of *morphomechanical fields*. This is a different approach from statistical spatial atlases frequently used to align information within a cohort [49, 50]. However, several fundamental questions and methodological obstacles remain unanswered. In particular, the sensitivity of the automatic classification of *morphomechanical fields* to intra-phenotypic variability, and its usefulness in establishing inter-phenotype differences need to be addressed in more detail. In particular, automating these analyses for cohorts of embryos requires systematic scanning across entire embryos to compensate for the different development rates of each embryo and its phenotype variability. Then, through the spatio-temporal registration of *fields* [54], it could be possible to infer robust phenotyping structures and assess the impact of dynamics variability into morphological configuration of tissues and organs.

UNDERSTANDING MULTI-SCALE EMBRYONIC DYNAMICS BY MACHINE LEARNING

Biological systems are often defined as networks of discrete elements or biochemical processes, which serve as a conceptual framework to glean mechanistic insight about their organization [128, 129]. Framing embryogenesis using this paradigm involves identifying morphogenetic events and fields [130], which can be diverse in nature, duration, and length-scale. Based on image data, one can define morphogenetic events as spatio-temporal spots of statistically abnormal behavior given a reference window. They may comprise subcellular or mesoscopic regions and a variable number of time frames and can be encapsulated by applying spatio-temporal connectivity [131]. When these fields are backtracked, they become unwound in time and space, allowing the discovery of intersections with past events and/or environmental cues. Likewise, forward tracking

of events can reveal cascade effects that propagate into one or more morphogenetic fields. The structured representation of digital embryos as spatio-temporally connected fields is a form of dimensionality reduction that fits machine learning-driven approaches.

Owing to recent advances in machine learning methods, computers can now perform human-like reasoning in tasks, such as conversation or gaming [132–134]. Deep learning (DL) architectures, such as Feed-Forward Networks, Recurrent Neural Networks (RNN) and Convolutional Neural Networks (CNN) can be applied to analyze interactions in the networks of elements employed in systems biology [135]. Consequently, the applications of DL to biomedicine are quickly becoming ubiquitous [136–139]. The analysis of high-throughput genomics data to study genetic variations in regulatory networks is no exception [140, 141]. A main barrier toward adopting DL in developmental biology has been its black-box nature, which does not easily permit inferring mechanisms or causal relationships, and makes it challenging to manipulate models to test hypotheses. Most ongoing efforts to mitigate these limitations can be encompassed under the epistemological paradigm of the Visible Neural Network [142]. One approach toward VNN is to represent the nodes in the network as “visible” data-driven models. This approach has been used to relate cell genotypes and phenotypes based on cell ontologies [143]. An alternate approach is to build the nodes in the network using explicit models based on theoretical or semi-empirical laws [144]. Both approaches allow for manipulating the inner machinery of the DL architecture, thereby facilitating hypothesis testing, the inference of causal relationships, and elucidating mechanisms. Furthermore, coupling DL model-driven architectures with multi-level structured training data can help reduce the amount of inputs, simplify the architecture and facilitate its interpretation [145]. Exhaustive simulations running on cloud technologies [146, 147] can leverage computational models and feed machine learning workflows to create multiple hypothesis to be tested *in-vivo*. In the case of embryo development, most theoretical and computational models are coarse grained and, thus, better suited to represent meso-scale and large-scale phenomena (see section Computational Models). Consequently, it could be beneficial to develop hybrid approaches in which cell-scale phenomena are modeled with DL. This type of bottom-up methodology has shown great potential to improve the prediction of chaotic deterministic systems, such as turbulent flow [148], but it should be noted that, epistemologically, it constitutes a transparent network of opaque nodes. Given that multiple relationships among genetic and biophysical processes evolve dynamically in space and time during morphogenesis, RNNs are a suitable approach to treat experimental data sequences. Several architectures of RNN have been proposed to improve training and solve the vanishing gradient problem through time [149]. LSTM comprise memory cells to infer long-term dependencies in sequences [150–152]. Gated Recurrent Units are another RNN architecture that addresses the long-term memory problem and outperforms LSTM in some applications [153, 154]. Sets of LSTM can be combined to design an encoder-decoder that

approaches the problem as a conversion of the input sequence into an intermediate fixed-length sequence (encoder) that can be further classified (decoder) [153, 155]. Recent advances in sequence analysis have been based on the idea of attention [156–161]. Attention architectures deal with long inputs by focusing on relevant frames of the sequence, eliminating the restriction of a fixed-length intermediate sequence, and leveraging intermediate states of the encoder as additional input to the decoder. Attention also provides clarity of the input-output relationships [156] and has shown promising results in end-to-end entailment of complex data sequences [162]. The transformer, an architecture without recurrence that relies on feed-forward layers and attention, has been proposed to exploit the potential of attention while allowing for massive parallelization [161, 163].

A key issue is how to pre-train [163, 164] and train these architectures with the data streams of morphogenesis. For instance, contextual bidirectional pre-training might facilitate characterizing strain-stress relationships given past and future tissue states (**Figures 2E,F**), in order to generate stress maps. Also, entailment of morphogenetic cues and mechanical events with fate map determination could be possible using the input defined by the profiles of cell trajectories, labeled according to a given morphogenetic field or a mature organ. In this regard, the scalability of biological domain tagging could introduce bottlenecks in the generation of training sets, particularly when considering the inherent variability of biological data. These tasks may require using several input vectors at the same time requiring extending current speech-oriented DL architectures [165, 166].

OUTLOOK

In this perspective, we have critically surveyed recent advances in computational methods for the characterization of embryogenesis, focusing on how to integrate data from biophysical measurements or models into cell fate maps. The ongoing surge in research efforts to quantify the biophysics of morphogenesis is leading to important methodological contributions and new insights about how genetics unfold into phenotypes. Despite these advances, the mechanistic description of morphogenesis remains challenging, given the dynamic and multi-scale nature of the underlying processes and the notable plasticity of immature cells. Thus, new methods are required to understand the interplay of physics, genetics, and epigenetics, leading to cell fate map determination. State-of-the-art imaging systems, image analyses, and computer vision methods are enabling the digital curation of multi-dimensional, high-resolution atlases of developing embryos. These data need to be structured in a systematic way to ensure experimental reproducibility and compatibility of different databases, which are necessary for statistically significant comparisons of large cohorts. In this sense, we posit that data analysis would benefit from a Lagrangian representation based on cell trajectories containing the cumulative histories of the spatio-temporal events and environmental cues cells experience along their

paths. This representation integrates spatial information into temporal sequences allowing for multi-scale discovery of *morphomechanical* fields.

Computational models offer a powerful toolbox to assimilate and explain experimental data, as well as to test new hypotheses. As these models grow in sophistication, they are beginning to predict and decipher the dynamics of developing embryos, based on multi-scale biophysical formalisms that can tackle spatio-temporal heterogeneity and complex mechanobiological interplays. These formalisms are benefitting from novel, minimally-invasive experimental approaches to measure the evolving mechanical properties of live embryos. However, the increasing diversity of models makes it difficult to identify, harmonize, and cross-validate a set of laws that govern the dynamics of morphogenesis. The lack of long-term maintenance of many open-source modeling codes makes this task additionally challenging.

In parallel, machine learning is quickly gaining traction as an alternative to classic model-driven computation to leverage intensive experimentation machine learning and causality inference tools [167, 168] can help test the completeness of models. In particular, these tools can elucidate morphomechanical domains formed by cells with similar dynamics, and link the formation of these domains with upstream biomechanical events. Deep learning (DL) architectures are becoming transparent and interpretable by nesting data-driven or model-driven visible nodes, and have been proven useful to discover causal links in other biological processes. For a holistic approach, DL is suitable to analyze spatio-temporal profiles, seek for events, discover patterns and identify dynamic entities. Multi-scale comparison of cohorts with model-driven DL architectures can be the basis to discover “missing data,” factors and critical spatio-temporal processes regulating phenotype configuration. Overall, the methodologies and approaches here discussed will have valuable practical applications for tissue engineering, stem cell research, genetics and behavior of diseases, drug studies, and synthetic life.

REFERENCES

- Keller R. Physical biology returns to morphogenesis. *Science*. (2012) 338:201–3. doi: 10.1126/science.1230718
- Davidson L, von Dassow M, Zhou J. Multi-scale mechanics from molecules to morphogenesis. *Int J Biochem Cell Biol*. (2009) 41:2147–62. doi: 10.1016/j.biocel.2009.04.015
- Blanchard GB, Adams RJ. Measuring the multi-scale integration of mechanical forces during morphogenesis. *Curr Opin Genet Dev*. (2011) 21:653–63. doi: 10.1016/j.gde.2011.08.008
- Gilmour D, Rembold M, Leptin M. From morphogen to morphogenesis and back. *Nature*. (2017) 541:311–20. doi: 10.1038/nature21348
- Heisenberg C-P, Bellaïche Y. Forces in tissue morphogenesis and patterning. *Cell*. (2013) 153:948–62. doi: 10.1016/j.cell.2013.05.008
- Lecuit T, Lenne P-F, Munro E. Force generation, transmission, and integration during cell and tissue morphogenesis. *Annu Rev Cell Dev Biol*. (2011) 27:157–84. doi: 10.1146/annurev-cellbio-100109-104027
- Heller E, Fuchs E. Tissue patterning and cellular mechanics. *J Cell Biol*. (2015) 211:219–31. doi: 10.1083/jcb.201506106
- Wozniak MA, Chen CS. Mechanotransduction in development: a growing role for contractility. *Nat Rev Mol Cell Biol*. (2009) 10:34–43. doi: 10.1038/nrm2592
- Mammoto T, Ingber DE. Mechanical control of tissue and organ development. *Development*. (2010) 137:1407–20. doi: 10.1242/dev.024166
- Zhang R, Han P, Yang H, Ouyang K, Lee D, Lin Y-F, et al. *In vivo* cardiac reprogramming contributes to zebrafish heart regeneration. *Nature*. (2013) 498:497–501. doi: 10.1038/nature12322
- Gálvez-Santisteban M, Chen D, Zhang R, Serrano R, Nguyen C, Zhao L, et al. Hemodynamic-mediated endocardial signaling controls *in vivo* myocardial reprogramming. *Elife*. (2019) 8:e44816. doi: 10.7554/eLife.44816
- Roca-Cusachs P, Conte V, Trepas X. Quantifying forces in cell biology. *Nat Cell Biol*. (2017) 19:742–51. doi: 10.1038/ncb3564
- Kwan KM, Otsuna H, Kidokoro H, Carney KR, Saijoh Y, Chien C-B. A complex choreography of cell movements shapes the vertebrate eye. *Development*. (2012) 139:359–72. doi: 10.1242/dev.071407
- Briggs JA, Weinreb C, Wagner DE, Megason S, Peshkin L, Kirschner MW, et al. The dynamics of gene expression in vertebrate embryogenesis at single-cell resolution. *Science*. (2018) 360:eaar5780. doi: 10.1126/science.aar5780

DATA AVAILABILITY STATEMENT

The dataset used for visualization in **Figure 1** and the free version of the visualization tool Mov-IT are freely available from the BioEmergences open workflow <http://bioemergences.iscpif.fr/bioemergences/openworkflow-index.php>. Data and tool are described in Faure et al. [36].

ETHICS STATEMENT

A dataset of a developing wild-type zebrafish embryo was presented in **Figure 1**. This dataset was produced by the BioEmergences lab (bioemergences.eu) as described in Faure et al. [36].

AUTHOR CONTRIBUTIONS

DP-E conceived the work, made the figures, and co-wrote the manuscript. JÁ advised for this work and co-wrote the manuscript.

FUNDING

This work was supported by NIH grants 1 R01 HD092216-01A1, NIH 1R01HL128630, 1R01HL130840, NIH 2R01 GM084227, and NSF grant NSF CBET – 1706436/1706571.

ACKNOWLEDGMENTS

We thank BioEmergences Lab-CNRS and Nadine Peyriéras for the joint work on computational developmental biology that inspired this work. We also thank Nicole Gorfinkiel for discussions on tissue mechanics. We thank the Biomedical Image Technologies Lab-UPM, Andres Santos, María Jesús Ledesma-Carbayo and Jose M. Goicolea for their collaboration on previous work.

15. Wagner DE, Weinreb C, Collins ZM, Briggs JA, Megason SG, Klein AM. Single-cell mapping of gene expression landscapes and lineage in the zebrafish embryo. *Science*. (2018) 360:981–7. doi: 10.1126/science.aar4362
16. Farrell JA, Wang Y, Riesenfeld SJ, Shekhar K, Regev A, and Schier AF. Single-cell reconstruction of developmental trajectories during zebrafish embryogenesis. *Science*. (2018) 360:eaar3131. doi: 10.1126/science.aar3131
17. Turing AM. The chemical basis of morphogenesis. *Bull Math Biol*. (1990) 52:153–97. doi: 10.1007/BF02459572
18. Kimmel CB, Warga RM, Schilling TF. Origin and organization of the zebrafish fate map. *Development*. (1990) 108:581–94.
19. Woo K, Fraser SE. Order and coherence in the fate map of the zebrafish nervous system. *Development*. (1995) 121:2595–609.
20. Sako K, Pradhan SJ, Barone V, Inglés-Prieto Á, Müller P, Ruprecht V, et al. Optogenetic control of nodal signaling reveals a temporal pattern of nodal signaling regulating cell fate specification during gastrulation. *Cell Rep*. (2016) 16:866–77. doi: 10.1016/j.celrep.2016.06.036
21. Chan CJ, Heisenberg C-P, Hiiragi T. Coordination of morphogenesis and cell-fate specification in development. *Curr Biol*. (2017) 27:R1024–35. doi: 10.1016/j.cub.2017.07.010
22. Keller PJ. Imaging morphogenesis: technological advances and biological insights. *Science*. (2013) 340:1234168. doi: 10.1126/science.1234168
23. Polacheck WJ, Chen CS. Measuring cell-generated forces: a guide to the available tools. *Nat Methods*. (2016) 13:415–23. doi: 10.1038/nmeth.3834
24. Villoutreix P, Delile J, Rizzi B, Duloquin L, Savy T, Bourguin P, et al. An integrated modelling framework from cells to organism based on a cohort of digital embryos. *Sci Rep*. (2016) 6:37438. doi: 10.1038/srep37438
25. Sharpe J. Computer modeling in developmental biology: growing today, essential tomorrow. *Development*. (2017) 144:4214–25. doi: 10.1242/dev.151274
26. Delile J, Herrmann M, Peyri eras N, Doursat R. A cell-based computational model of early embryogenesis coupling mechanical behaviour and gene regulation. *Nat Commun*. (2017) 8:13929. doi: 10.1038/ncomms13929
27. Yeh Y-T, Serrano R, Fran ois J, Chiu J-J, Li Y-SJ, Del Álamo JC, et al. Three-dimensional forces exerted by leukocytes and vascular endothelial cells dynamically facilitate diapedesis. *Proc Natl Acad Sci USA*. (2018) 115:133–8. doi: 10.1073/pnas.1717489115
28. Latorre E, Kale S, Casares L, G omez-Gonz alez M, Uroz M, Valon L, et al. Active superelasticity in three-dimensional epithelia of controlled shape. *Nature*. (2018) 563:203–8. doi: 10.1038/s41586-018-0671-4
29. Forgacs G, Foty RA, Shafir Y, Steinberg MS. Viscoelastic properties of living embryonic tissues: a quantitative study. *Biophys J*. (1998) 74:2227–34. doi: 10.1016/S0006-3495(98)77932-9
30. Marmottant P, Mgharbel A, K fer J, Audren B, Rieu J-P, Vial J-C, et al. The role of fluctuations and stress on the effective viscosity of cell aggregates. *Proc Natl Acad Sci USA*. (2009) 106:17271–5. doi: 10.1073/pnas.0902085106
31. Wu P-H, Aroush DR-B, Asnacios A, Chen W-C, Dokukin ME, Doss BL, et al. A comparison of methods to assess cell mechanical properties. *Nat Methods*. (2018) 15:491–8. doi: 10.1038/s41592-018-0015-1
32. Olivier N, Luengo-Oroz MA, Duloquin L, Faure E, Savy T, Veilleux I, et al. Cell lineage reconstruction of early zebrafish embryos using label-free nonlinear microscopy. *Science*. (2010) 329:967–71. doi: 10.1126/science.1189428
33. Supatto W, Truong TV, D barre D, Beaupaire E. Advances in multiphoton microscopy for imaging embryos. *Curr Opin Genet Dev*. (2011) 21:538–48. doi: 10.1016/j.gde.2011.08.003
34. Gao L, Shao L, Chen B-C, Betzig E. 3D live fluorescence imaging of cellular dynamics using Bessel beam plane illumination microscopy. *Nat Protoc*. (2014) 9:1083–1101. doi: 10.1038/nprot.2014.087
35. Wolff C, Tinevez J-Y, Pietzsch T, Stamataki E, Harich B, Guignard L, et al. Multi-view light-sheet imaging and tracking with the MaMuT software reveals the cell lineage of a direct developing arthropod limb. *Elife*. (2018) 7:e34410. doi: 10.7554/eLife.34410
36. Faure E, Savy T, Rizzi B, Melani C, Sta ov a O, Fabr ges D, et al. A workflow to process 3D+ time microscopy images of developing organisms and reconstruct their cell lineage. *Nat Commun*. (2016) 7:8674. doi: 10.1038/ncomms9674
37. Amat F, Lemon W, Mossing DP, McDole K, Wan Y, Branson K, et al. Fast, accurate reconstruction of cell lineages from large-scale fluorescence microscopy data. *Nat Methods*. (2014) 11:951–8. doi: 10.1038/nmeth.3036
38. Tinevez J-Y, Pietzsch T, Rueden C. *MaMuT*. Github (2018). Available online at: <https://github.com/fiji/MaMuT> (accessed September 1, 2019).
39. Stegmaier J, Amat F, Lemon WC, McDole K, Wan Y, Teodoro G, et al. Real-time three-dimensional cell segmentation in large-scale microscopy data of developing embryos. *Dev Cell*. (2016) 36:225–40. doi: 10.1016/j.devcel.2015.12.028
40. Ulman V, Ma ka M, Magnusson KE, Ronneberger O, Haubold C, Harder N, et al. An objective comparison of cell-tracking algorithms. *Nat Methods*. (2017) 14:1141–52. doi: 10.1038/nmeth.4473
41. Dufour AC, Jonker AH, Olivo-Marin J-C. Deciphering tissue morphodynamics using bioimage informatics. *Philos Trans R Soc B Biol Sci*. (2017) 372:20150512. doi: 10.1098/rstb.2015.0512
42. Schott B, Traub M, Schlagenhauf C, Takamiya M, Anritter T, Bartschat A, et al. EmbryoMiner: a new framework for interactive knowledge discovery in large-scale cell tracking data of developing embryos. *PLoS Comput Biol*. (2018) 14:e1006128. doi: 10.1371/journal.pcbi.1006128
43. Leggio B, Laussu J, Carlier A, Godin C, Lemaire P, Faure E. MorphoNet: an interactive online morphological browser to explore complex multi-scale data. *Nat Commun*. (2019) 10:2812. doi: 10.1038/s41467-019-10668-1
44. Tomer R, Khairy K, Amat F, Keller PJ. Quantitative high-speed imaging of entire developing embryos with simultaneous multiview light-sheet microscopy. *Nat Methods*. (2012) 9:755–63. doi: 10.1038/nmeth.2062
45. Wu Y, Chandris P, Winter PW, Kim EY, Jaumouill  V, Kumar A, et al. Simultaneous multiview capture and fusion improves spatial resolution in wide-field and light-sheet microscopy. *Optica*. (2016) 3:897–910. doi: 10.1364/OPTICA.3.000897
46. Rubio-Guivernau JL, Gurchenkov V, Luengo-Oroz MA, Duloquin L, Bourguin P, Santos A, et al. Wavelet-based image fusion in multi-view three-dimensional microscopy. *Bioinformatics*. (2011) 28:238–45. doi: 10.1093/bioinformatics/btr609
47. Schmier C, Stamataki E, Tomancak P. Open-source solutions for SPIMage processing. *Methods Cell Biol*. (2014) 505–29. doi: 10.1016/B978-0-12-420138-5.00027-6
48. Liu T-L, Upadhyayula S, Milkie DE, Singh V, Wang K, Swinburne IA, et al. Observing the cell in its native state: imaging subcellular dynamics in multicellular organisms. *Science*. (2018) 360:eaq1392. doi: 10.1126/science.aq1392
49. Ronneberger O, Liu K, Rath M, Rue  D, Mueller T, Skibbe H, et al. ViBE-Z: a framework for 3D virtual colocalization analysis in zebrafish larval brains. *Nat Methods*. (2012) 9:735–42. doi: 10.1038/nmeth.2076
50. Castro-Gonz alez C, Luengo-Oroz MA, Duloquin L, Savy T, Rizzi B, Desnoul ez S, et al. A digital framework to build, visualize and analyze a gene expression atlas with cellular resolution in zebrafish early embryogenesis. *PLoS Comput Biol*. (2014) 10:e1003670. doi: 10.1371/journal.pcbi.1003670
51. Power RM, Huiskens J. A guide to light-sheet fluorescence microscopy for multiscale imaging. *Nat Methods*. (2017) 14:360–73. doi: 10.1038/nmeth.4224
52. Ove ka M, von Wangenheim D, Toman ak P,  amajov a O, Komis G,  amaj J. Multiscale imaging of plant development by light-sheet fluorescence microscopy. *Nature Plants*. (2018) 4:639–50. doi: 10.1038/s41477-018-0238-2
53. Pastor Escuredo D. *Methods for the Analysis of Multi-Scale Cell Dynamics From Fluorescence Microscopy Images*. Madrid: Universidad Polit cnica de Madrid, ETSIT (2015).
54. Ledesma-Carbayo MJ, Kybic J, Desco M, Santos A, Suhling M, Hunziker P, et al. Spatio-temporal nonrigid registration for ultrasound cardiac motion estimation. *IEEE Trans Med Imaging*. (2005) 24:1113–26. doi: 10.1109/TMI.2005.852050
55. Blanchard GB, Kabla AJ, Schultz NL, Butler LC, Sanson B, Gorfinkiel N, et al. Tissue tectonics: morphogenetic strain rates, cell shape change and intercalation. *Nat Methods*. (2009) 6:458–64. doi: 10.1038/nmeth.1327
56. He B, Doubrovinski K, Polyakov O, Wieschaus E. Apical constriction drives tissue-scale hydrodynamic flow to mediate cell elongation. *Nature*. (2014) 508:392–6. doi: 10.1038/nature13070

57. Pastor-Escuredo D, Lombardot B, Savy T, Boyreau A, Goicolea JM, Santos A, et al. Kinematic analysis of cell lineage reveals coherent and robust mechanical deformation patterns in zebrafish gastrulation. *bioRxiv* 054353. doi: 10.1101/054353
58. Stegmaier J, Spina TV, Falcão AX, Bartschat A, Mikut R, Meyerowitz E, et al. Cell segmentation in 3D confocal images using supervoxel merge-forests with CNN-based hypothesis selection. In: *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)*. Washington, DC: IEEE (2018). p. 382–386.
59. Caicedo JC, Roth J, Goodman A, Becker T, Karhohs KW, Broisin M, et al. Evaluation of deep learning strategies for nucleus segmentation in fluorescence images. *J Quant Cell Sci*. 95:952–65. doi: 10.1002/cyto.a.23863
60. Etournay R, Merkel M, Popović M, Brandl H, Dye NA, Aigouy B, et al. TissueMiner: a multiscale analysis toolkit to quantify how cellular processes create tissue dynamics. *Elife*. (2016) 5:e14334. doi: 10.7554/eLife.14334.033
61. Linares M, Postigo M, Cuadrado D, Ortiz-Ruiz A, Gil-Casanova S, Vladimirov A, et al. Collaborative intelligence and gamification for on-line malaria species differentiation. *Malar J*. (2019) 18:21. doi: 10.1186/s12936-019-2662-9
62. Sugimura K, Lenne P-F, Graner F. Measuring forces and stresses *in situ* in living tissues. *Development*. (2016) 143:186–96. doi: 10.1242/dev.119776
63. Rauzi M, Verant P, Lecuit T, Lenne P-F. Nature and anisotropy of cortical forces orienting *Drosophila* tissue morphogenesis. *Nat Cell Biol*. (2008) 10:1401–10. doi: 10.1038/ncb1798
64. Hutson MS, Veldhuis J, Ma X, Lynch HE, Cranston PG, Brodland GW. Combining laser microsurgery and finite element modeling to assess cell-level epithelial mechanics. *Biophys J*. (2009) 97:3075–85. doi: 10.1016/j.bpj.2009.09.034
65. Behrndt M, Salbreux G, Campinho P, Hauschild R, Oswald F, Roensch J, et al. Forces driving epithelial spreading in zebrafish gastrulation. *Science*. (2012) 338:257–60. doi: 10.1126/science.1224143
66. Levayer R, Lecuit T. Biomechanical regulation of contractility: spatial control and dynamics. *Trends Cell Biol*. (2012) 22:61–81. doi: 10.1016/j.tcb.2011.10.001
67. Machado PF, Duque J, Étienne J, Martinez-Arias A, Blanchard GB, Gorfinkiel N. Emergent material properties of developing epithelial tissues. *BMC Biol*. (2015) 13:98. doi: 10.1186/s12915-015-0200-y
68. Jurchenko C, Salaita KS. Lighting up the force: investigating mechanisms of mechanotransduction using fluorescent tension probes. *Mol Cell Biol*. (2015) 35:2570–82. doi: 10.1128/MCB.00195-15
69. Dembo M, Wang Y-L. Stresses at the cell-to-substrate interface during locomotion of fibroblasts. *Biophys J*. (1999) 76:2307–16. doi: 10.1016/S0006-3495(99)77386-8
70. Butler JP, Tolic-Nørrelykke IM, Fabry B, Fredberg JJ. Traction fields, moments, and strain energy that cells exert on their surroundings. *Am J Physiol Cell Physiol*. (2002) 282:C595–605. doi: 10.1152/ajpcell.00270.2001
71. Del Álamo JC, Meili R, Alonso-Latorre B, Rodríguez-Rodríguez J, Aliseda A, Firtel RA, et al. Spatio-temporal analysis of eukaryotic cell motility by improved force cytometry. *Proc Natl Acad Sci*. (2007) 104:13343–8. doi: 10.1073/pnas.0705815104
72. Trepast X, Wasserman MR, Angelini TE, Millet E, Weitz DA, Butler JP, et al. Physical forces during collective cell migration. *Nat Phys*. (2009) 5:426–30. doi: 10.1038/nphys1269
73. Brodland GW, Conte V, Cranston PG, Veldhuis J, Narasimhan S, Hutson MS, et al. Video force microscopy reveals the mechanics of ventral furrow invagination in *Drosophila*. *Proc Natl Acad Sci USA*. (2010) 107:22111–6. doi: 10.1073/pnas.1006591107
74. Hur SS, Zhao Y, Li Y-S, Botvinick E, Chien S. Live cells exert 3-dimensional traction forces on their substrata. *Cell Mol Bioeng*. (2009) 2:425–36. doi: 10.1007/s12195-009-0082-6
75. Serrano R, Aung A, Varghese S, del Álamo JC. Three-dimensional monolayer stress microscopy. *Biophys J*. (2016) 110:330a. doi: 10.1016/j.bpj.2015.11.1776
76. Heller D, Hoppe A, Restrepo S, Gatti L, Tournier AL, Tapon N, et al. EpiTools: an open-source image analysis toolkit for quantifying epithelial growth dynamics. *Dev Cell*. (2016) 36:103–16. doi: 10.1016/j.devcel.2015.12.012
77. Bosveld F, Bonnet I, Guirao B, Tlili S, Wang Z, Petalot A, et al. Mechanical control of morphogenesis by Fat/Dachsous/Four-jointed planar cell polarity pathway. *Science*. (2012) 336:724–7. doi: 10.1126/science.1221071
78. Blanchard GB. Taking the strain: quantifying the contributions of all cell behaviours to changes in epithelial shape. *Philos Trans R Soc B Biol Sci*. (2017) 372:20150513. doi: 10.1098/rstb.2015.0513
79. Tetley RJ, Blanchard GB, Fletcher AG, Adams RJ, Sanson B. Unipolar distributions of junctional Myosin II identify cell stripe boundaries that drive cell intercalation throughout *Drosophila* axis extension. *Elife*. (2016) 5:e12094. doi: 10.7554/eLife.12094
80. Tlili S, Gay C, Graner F, Marcq P, Molino F, Saramito P. Colloquium: mechanical formalisms for tissue dynamics. *Eur Phys J E*. (2015) 38:33. doi: 10.1140/epje/i2015-15033-4
81. Chiou KK, Hufnagel L, Shraiman BI. Mechanical stress inference for two dimensional cell arrays. *PLoS Comput Biol*. (2012) 8:e1002512. doi: 10.1371/journal.pcbi.1002512
82. Ishihara S, Sugimura K. Bayesian inference of force dynamics during morphogenesis. *J Theor Biol*. (2012) 313:201–11. doi: 10.1016/j.jtbi.2012.08.017
83. Ishihara S, Sugimura K, Cox S, Bonnet I, Bellaiche Y, Graner F. Comparative study of non-invasive force and stress inference methods in tissue. *Eur Phys J E*. (2013) 36:45. doi: 10.1140/epje/i2013-13045-8
84. Brodland GW, Veldhuis JH, Kim S, Perrone M, Mashburn D, Hutson MS. CellFIT: a cellular force-inference toolkit using curvilinear cell boundaries. *PLoS ONE*. (2014) 9:e99116. doi: 10.1371/journal.pone.0099116
85. Guirao B, Rigaud SU, Bosveld F, Bailles A, Lopez-Gay J, Ishihara S, et al. Unified quantitative characterization of epithelial tissue development. *Elife*. (2015) 4:e08519. doi: 10.7554/eLife.08519
86. Veldhuis JH, Ehsandar A, Maître J-L, Hiiragi T, Cox S, Brodland GW. Inferring cellular forces from image stacks. *Philos Trans R Soc B Biol Sci*. (2017) 372:20160261. doi: 10.1098/rstb.2016.0261
87. Mongera A, Rowghanian P, Gustafson HJ, Shelton E, Kealhofer DA, Carn EK, et al. A fluid-to-solid jamming transition underlies vertebrate body axis elongation. *Nature*. (2018) 561:401–5. doi: 10.1038/s41586-018-0479-2
88. Campàs O, Mammoto T, Hasso S, Sperling RA, O'Connell D, Bischof AG, et al. Quantifying cell-generated mechanical forces within living embryonic tissues. *Nat Methods*. (2014) 11:183–9. doi: 10.1038/nmeth.2761
89. Dolega M, Delarue M, Ingremeau F, Prost J, Delon A, Cappello G. Cell-like pressure sensors reveal increase of mechanical stress towards the core of multicellular spheroids under compression. *Nat Commun*. (2017) 8:14056. doi: 10.1038/ncomms14056
90. Serwane F, Mongera A, Rowghanian P, Kealhofer DA, Lucio AA, Hockenberg ZM, et al. *In vivo* quantification of spatially varying mechanical properties in developing tissues. *Nat Methods*. (2017) 14:181–6. doi: 10.1038/nmeth.4101
91. Byrne H, Drasdo D. Individual-based and continuum models of growing cell populations: a comparison. *J Math Biol*. (2009) 58:657. doi: 10.1007/s00285-008-0212-0
92. Van Liedekerke, P. (2019). *Quantitative Modeling of Cell and Tissue Mechanics With Agent-Based Models*. Paris: Inria Paris, Sorbonne Université.
93. Van Liedekerke P, Palm M, Jagiella N, Drasdo D. Simulating tissue mechanics with agent-based models: concepts, perspectives and some novel results. *Comput Part Mech*. (2015) 2:401–44. doi: 10.1007/s40571-015-0082-3
94. Drasdo D, Höhme S. A single-cell-based model of tumor growth *in vitro*: monolayers and spheroids. *Phys Biol*. (2005) 2:133–47. doi: 10.1088/1478-3975/2/3/001
95. Van Liedekerke P, Neitsch J, Johann T, Warnt E, Grosser S, Valverde IG, et al. Quantifying the mechanics and growth of cells and tissues in 3D using high resolution computational models. *bioRxiv* 470559. doi: 10.1101/470559
96. Swat MH, Thomas GL, Belmonte JM, Shirinifard A, Hmeljak D, Glazier JA. Multi-scale modeling of tissues using CompuCell3D. *Methods Cell Biol*. (2012) 325–66. doi: 10.1016/B978-0-12-388403-9.00013-8
97. Voss-Böhme A. Multi-scale modeling in morphogenesis: a critical analysis of the cellular Potts model. *PLoS ONE*. (2012) 7:e42852. doi: 10.1371/journal.pone.0042852

98. Graner F, Glazier JA. Simulation of biological cell sorting using a two-dimensional extended Potts model. *Phys Rev Lett.* (1992) 69:2013. doi: 10.1103/PhysRevLett.69.2013
99. Käfer J, Hayashi T, Marée AF, Carthew RW, Graner F. Cell adhesion and cortex contractility determine cell patterning in the *Drosophila* retina. *Proc Natl Acad Sci USA.* (2007) 104:18549–54. doi: 10.1073/pnas.0704235104
100. Krieg M, Arboleda-Estudillo Y, Puech P-H, Käfer J, Graner F, Müller D, et al. Tensile forces govern germ-layer organization in zebrafish. *Nat Cell Biol.* (2008) 10:429–36. doi: 10.1038/ncb1705
101. Fletcher AG, Cooper F, Baker RE. Mechanocellular models of epithelial morphogenesis. *Philos Trans R Soc B Biol Sci.* (2017) 372:20150519. doi: 10.1098/rstb.2015.0519
102. Alt S, Ganguly P, Salbreux G. Vertex models: from cell mechanics to tissue morphogenesis. *Philos Trans R Soc B Biol Sci.* (2017) 372:20150520. doi: 10.1098/rstb.2015.0520
103. Fletcher AG, Osterfield M, Baker RE, Shvartsman SY. Vertex models of epithelial morphogenesis. *Biophys J.* (2014) 106:2291–304. doi: 10.1016/j.bpj.2013.11.4498
104. Nagai T, Honda H. A dynamic cell model for the formation of epithelial tissues. *Philos Mag B.* (2001) 81:699–719. doi: 10.1080/13642810108205772
105. Farhadifar R, Röper J-C, Aigouy B, Eaton S, Jülicher F. The influence of cell mechanics, cell-cell interactions, and proliferation on epithelial packing. *Curr Biol.* (2007) 17:2095–104. doi: 10.1016/j.cub.2007.11.049
106. Ishimoto Y, Morishita Y. Bubbly vertex dynamics: a dynamical and geometrical model for epithelial tissues with curved cell shapes. *Phys Rev E.* (2014) 90:052711. doi: 10.1103/PhysRevE.90.052711
107. Perrone MC, Veldhuis JH, Brodland GW. Non-straight cell edges are important to invasion and engulfment as demonstrated by cell mechanics model. *Biomech Model Mechanobiol.* (2016) 15:405–18. doi: 10.1007/s10237-015-0697-6
108. Tanaka S, Sichau D, Iber D. LBIBCell: a cell-based simulation environment for morphogenetic problems. *Bioinformatics.* (2015) 31:2340–7. doi: 10.1093/bioinformatics/btv147
109. Brezavšek AH, Rauzi M, Leptin M, Zihler P. A model of epithelial invagination driven by collective mechanics of identical cells. *Biophys J.* (2012) 103:1069–77. doi: 10.1016/j.bpj.2012.07.018
110. Okuda S, Inoue Y, Eiraku M, Sasai Y, Adachi T. Apical contractility in growing epithelium supports robust maintenance of smooth curvatures against cell-division-induced mechanical disturbance. *J Biomech.* (2013) 46:1705–13. doi: 10.1016/j.jbiomech.2013.03.035
111. Murisic N, Hakim V, Kevrekidis IG, Shvartsman SY, Audoly B. From discrete to continuum models of three-dimensional deformations in epithelial sheets. *Biophys J.* (2015) 109:154–63. doi: 10.1016/j.bpj.2015.05.019
112. Monier B, Gettings M, Gay G, Mangeat T, Schott S, Guarner A, et al. Apico-basal forces exerted by apoptotic cells drive epithelium folding. *Nature.* (2015) 518:245–8. doi: 10.1038/nature14152
113. Rauzi M, Krzic U, Saunders TE, Krajnc M, Zihler P, Hufnagel L, et al. Embryo-scale tissue mechanics during *Drosophila* gastrulation movements. *Nat Commun.* (2015) 6:8677. doi: 10.1038/ncomms9677
114. Misra M, Audoly B, Kevrekidis IG, Shvartsman SY. Shape transformations of epithelial shells. *Biophys J.* (2016) 110:1670–8. doi: 10.1016/j.bpj.2016.03.009
115. Hufnagel L, Teleman AA, Rouault H, Cohen SM, Shraiman BI. On the mechanism of wing size determination in fly development. *Proc Natl Acad Sci USA.* (2007) 104:3835–40. doi: 10.1073/pnas.0607134104
116. Aegerter-Wilmsen T, Smith AC, Christen AJ, Aegerter CM, Hafen E, Basler K. Exploring the effects of mechanical feedback on epithelial topology. *Development.* (2010) 137:499–506. doi: 10.1242/dev.041731
117. Li Y, Naveed H, Kachalo S, Xu LX, Liang J. Mechanisms of regulating cell topology in proliferating epithelia: impact of division plane, mechanical forces, and cell memory. *PLoS ONE.* (2012) 7:e43108. doi: 10.1371/journal.pone.0043108
118. Marinari E, Mehonic A, Curran S, Gale J, Duke T, Baum B. Live-cell delamination counterbalances epithelial growth to limit tissue overcrowding. *Nature.* (2012) 484:542–45. doi: 10.1038/nature10984
119. Aigouy B, Farhadifar R, Staple DB, Sagner A, Röper J-C, Jülicher F, et al. Cell flow reorients the axis of planar polarity in the wing epithelium of *Drosophila*. *Cell.* (2010) 142:773–86. doi: 10.1016/j.cell.2010.07.042
120. Landsberg KP, Farhadifar R, Ranft J, Umetsu D, Widmann TJ, Bittig T, et al. Increased cell bond tension governs cell sorting at the *Drosophila* anteroposterior compartment boundary. *Curr Biol.* (2009) 19:1950–5. doi: 10.1016/j.cub.2009.10.021
121. Zhao J, Cao Y, DiPietro LA, Liang J. Dynamic cellular finite-element method for modelling large-scale cell migration and proliferation under the control of mechanical and biochemical cues: a study of re-epithelialization. *J R Soc Interface.* (2017) 14. doi: 10.1098/rsif.2016.0959
122. Marchetti MC, Joanny J-F, Ramaswamy S, Liverpool TB, Prost J, Rao M, et al. Hydrodynamics of soft active matter. *Rev Mod Phys.* (2013) 85:1143. doi: 10.1103/RevModPhys.85.1143
123. Blanchard GB, Fletcher AG, Schumacher LJ. The devil is in the mesoscale: mechanical and behavioural heterogeneity in collective cell movement. *Semin Cell Dev Biol.* (2018) 93:46–54. doi: 10.1016/j.semcdb.2018.06.003
124. Conte V, Muñoz JJ, Miodownik M. A 3D finite element model of ventral furrow invagination in the *Drosophila melanogaster* embryo. *J Mech Behav Biomed Mater.* (2008) 1:188–98. doi: 10.1016/j.jmbbm.2007.10.002
125. González-Valverde I, García-Aznar JM. Mechanical modeling of collective cell migration: an agent-based and continuum material approach. *Comput Methods Appl Mech Eng.* (2018) 337:246–62. doi: 10.1016/j.cma.2018.03.036
126. Shadden SC, Lekien F, Marsden JE. Definition and properties of Lagrangian coherent structures from finite-time Lyapunov exponents in two-dimensional aperiodic flows. *Physica D.* (2005) 212:271–304. doi: 10.1016/j.physd.2005.10.007
127. Haller G. Lagrangian coherent structures. *Annu Rev Fluid Mech.* (2015) 47:137–62. doi: 10.1146/annurev-fluid-010313-141322
128. Kitano H. Systems biology: a brief overview. *Science.* (2002) 295:1662–4. doi: 10.1126/science.1069492
129. Biasuz K, Leggio B, Faure E, Lemaire P. The “computable egg”: myth or useful concept? *Curr Opin Syst Biol.* (2018) 11:91–7. doi: 10.1016/j.coisb.2018.09.003
130. Vecchi D, Hernández I. The epistemological resilience of the concept of morphogenetic field. In: *Towards a Theory of Development*. Oxford, UK: Oxford University Press (2014). p. 79.
131. Luengo-Oroz MA, Pastor-Escuredo D, Castro-Gonzalez C, Faure E, Savy T, Lombardot B, et al. \$3D+t\$ morphological processing: applications to embryogenesis image analysis. *IEEE Trans Image Proc.* (2012) 21:3518–30. doi: 10.1109/TIP.2012.2197007
132. Hinton G, Deng L, Yu D, Dahl G, Mohamed AR, Jaitly N, et al. Deep neural networks for acoustic modeling in speech recognition. *IEEE Signal Proc Mag.* (2012) 29:82–97. doi: 10.1109/MSP.2012.2205597
133. Sainath TN, Mohamed AR, Kingsbury B, Ramabhadran B. Deep convolutional neural networks for LVCSR. In: *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*. Vancouver, BC: IEEE (2013). p. 8614–8.
134. Silver D, Huang A, Maddison CJ, Guez A, Sifre L, Van Den Driessche G, et al. Mastering the game of Go with deep neural networks and tree search. *Nature.* (2016) 529:484–9. doi: 10.1038/nature16961
135. LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature.* (2015) 521:436–44. doi: 10.1038/nature14539
136. Mamoshina P, Vieira A, Putin E, Zhavoronkov A. Applications of deep learning in biomedicine. *Mol Pharm.* (2016) 13:1445–54. doi: 10.1021/acs.molpharmaceut.5b00982
137. Min S, Lee B, Yoon S. Deep learning in bioinformatics. *Brief Bioinformatics.* (2017) 18:851–69. doi: 10.1093/bib/bbw068
138. Wainberg M, Merico D, Delong A, Frey BJ. Deep learning in biomedicine. *Nat Biotechnol.* (2018) 36:829–38. doi: 10.1038/nbt.4233
139. Zitnik M, Nguyen F, Wang B, Leskovec J, Goldenberg A, Hoffman MM. Machine learning for integrating data in biology and medicine: Principles, practice, and opportunities. *Inform Fus.* (2019) 50:71–91. doi: 10.1016/j.inffus.2018.09.012
140. Angermueller C, Pärnamaa T, Parts L, Stegle O. Deep learning for computational biology. *Mol Syst Biol.* (2016) 12:878. doi: 10.15252/msb.20156651
141. Eraslan G, Avsec Ž, Gagnew J, Theis FJ. Deep learning: new computational modelling techniques for genomics. *Nat Rev Genet.* 20:389–403. doi: 10.1038/s41576-019-0122-6

142. Michael KY, Ma J, Fisher J, Kreisberg JF, Raphael BJ, Ideker T. Visible machine learning for biomedicine. *Cell*. (2018) 173:1562–5. doi: 10.1016/j.cell.2018.05.056
143. Ma J, Yu MK, Fong S, Ono K, Sage E, Demchak B, et al. Using deep learning to model the hierarchical structure and function of a cell. *Nat Methods*. (2018) 15:290–8. doi: 10.1038/nmeth.4627
144. Gazestani VH, Lewis NE. From genotype to phenotype: Augmenting deep learning with networks and systems biology. *Curr Opin Syst Biol*. (2019) 15:68–73. doi: 10.1016/j.coisb.2019.04.001
145. Wang D, Liu S, Warrell J, Won H, Shi X, Navarro FC, et al. (2018). Comprehensive functional genomic resource and integrative model for the human brain. *Science*. 362:eaat8464. doi: 10.1126/science.aat8464
146. Coulier A, Hellander A. Orchestral: a lightweight framework for parallel simulations of cell-cell communication. In: *2018 IEEE 14th International Conference on e-Science (e-Science)*. Amsterdam: IEEE (2018), 168–176.
147. Ghaffarizadeh A, Heiland R, Friedman SH, Mumenthaler SM, Macklin P. PhysiCell: an open source physics-based cell simulator for 3-D multicellular systems. *PLoS Comput Biol*. (2018) 14:e1005991. doi: 10.1371/journal.pcbi.1005991
148. Tracey BD, Duraisamy K, and Alonso JJ. A machine learning strategy to assist turbulence model development. In: *53rd AIAA Aerospace Sciences Meeting*. Kissimmee, FL (2015). p. 1287.
149. Jozefowicz R, Zaremba W, Sutskever I. An empirical exploration of recurrent network architectures. In: *International Conference on Machine Learning*. Lille (2015). p. 2342–50.
150. Hochreiter S, Schmidhuber J. Long short-term memory. *Neural Comput*. (1997) 9:1735–80. doi: 10.1162/neco.1997.9.8.1735
151. Sak H, Senior A, Beaufays F. Long short-term memory recurrent neural network architectures for large scale acoustic modeling. In: *Fifteenth Annual Conference of the International Speech Communication Association*. Singapore (2014).
152. Malhotra P, Vig L, Shroff G, and Agarwal, P. Long short term memory networks for anomaly detection in time series. In: *Proceedings: Presses universitaires de Louvain*. Louvain (2015). p. 89.
153. Cho K, Van Merriënboer B, Gulcehre C, Bahdanau D, Bougares F, Schwenk H, et al. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv [Preprint] arXiv:1406.1078* (2014). doi: 10.3115/v1/D14-1179
154. Chung J, Gulcehre C, Cho K, and Bengio Y. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv [Preprint] arXiv:1412.3555* (2014).
155. Sutskever I, Vinyals O, Le QV. Sequence to sequence learning with neural networks. In: *Advances in Neural Information Processing Systems*. Montréal, QC (2014). p. 3104–12.
156. Bahdanau D, Cho K, and Bengio Y. Neural machine translation by jointly learning to align and translate. *arXiv [Preprint] arXiv:1409.0473* (2014).
157. Xu K, Ba J, Kiros R, Cho K, Courville A, Salakhudinov R, et al. Show, attend and tell: neural image caption generation with visual attention. In: *International Conference on Machine Learning*. Lille (2015). p. 2048–57.
158. Hermann KM, Kocisky T, Grefenstette E, Espeholt L, Kay W, Suleyman M, et al. Teaching machines to read and comprehend. In: *Advances in Neural Information Processing Systems*. Montréal, QC (2015). p. 1693–1701.
159. Luong M T, Pham H, Manning CD. Effective approaches to attention-based neural machine translation. *arXiv [Preprint] arXiv:1508.04025* (2015). doi: 10.18653/v1/D15-1166
160. You Q, Jin H, Wang Z, Fang C, and Luo, J. Image captioning with semantic attention. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Las Vegas, NV (2016). p. 4651–9.
161. Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, et al. Attention is all you need. In: *Advances in Neural Information Processing Systems*. Long Beach, CA (2017). p. 5998–6008.
162. Rocktäschel T, Grefenstette E, Hermann KM, Kočický T, Blunsom, P. Reasoning about entailment with neural attention. *arXiv [Preprint] arXiv:1509.06664* (2015).
163. Devlin J, Chang MW, Lee K, Toutanova K. Bert: pre-training of deep bidirectional transformers for language understanding. *arXiv [Preprint] arXiv:1810.04805* (2018).
164. Dai AM, Le QV. Semi-supervised sequence learning. In: *Advances in Neural Information Processing Systems*. Montréal, QC (2015). p. 3079–87.
165. Chorowski JK, Bahdanau D, Serdyuk D, Cho K, Bengio Y. Attention-based models for speech recognition. In: *Advances in Neural Information Processing Systems*. Montréal, QC (2015). p. 577–85.
166. Rush AM, Chopra S, Weston J. A neural attention model for abstractive sentence summarization. *arXiv [Preprint] arXiv:1509.00685* (2015). doi: 10.18653/v1/D15-1044
167. Granger CW. Causality, cointegration, and control. *J Econ Dyn Control*. (1988) 12:551–9. doi: 10.1016/0165-1889(88)90055-3
168. Pearl J. *Causality: Models, Reasoning and Inference*. Cambridge, MA: Cambridge University Press (2000).

Conflict of Interest: The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2020 Pastor-Escuredo and del Álamo. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.



Unsupervised Learning Facilitates Neural Coordination Across the Functional Clusters of the *C. elegans* Connectome

Alejandro Morales^{1,2*} and Tom Froese¹

¹ Embodied Cognitive Science Unit, Okinawa Institute of Science and Technology Graduate University, Okinawa, Japan,

² Computer Science and Engineering Postgraduate Program, National Autonomous University of Mexico, Mexico City, Mexico

OPEN ACCESS

Edited by:

Georg Martius,
Max Planck Institute for Intelligent
Systems, Germany

Reviewed by:

Pablo Varona,
Autonomous University of
Madrid, Spain
Sam Neymotin,
Nathan Kline Institute for Psychiatric
Research, United States

*Correspondence:

Alejandro Morales
alejandroe@ciencias.unam.mx

Specialty section:

This article was submitted to
Computational Intelligence in
Robotics,
a section of the journal
Frontiers in Robotics and AI

Received: 06 December 2019

Accepted: 09 March 2020

Published: 02 April 2020

Citation:

Morales A and Froese T (2020)
Unsupervised Learning Facilitates
Neural Coordination Across the
Functional Clusters of the *C. elegans*
Connectome. *Front. Robot. AI* 7:40.
doi: 10.3389/frobt.2020.00040

Modeling of complex adaptive systems has revealed a still poorly understood benefit of unsupervised learning: when neural networks are enabled to form an associative memory of a large set of their own attractor configurations, they begin to reorganize their connectivity in a direction that minimizes the coordination constraints posed by the initial network architecture. This self-optimization process has been replicated in various neural network formalisms, but it is still unclear whether it can be applied to biologically more realistic network topologies and scaled up to larger networks. Here we continue our efforts to respond to these challenges by demonstrating the process on the connectome of the widely studied nematode worm *C. elegans*. We extend our previous work by considering the contributions made by hierarchical partitions of the connectome that form functional clusters, and we explore possible beneficial effects of inter-cluster inhibitory connections. We conclude that the self-optimization process can be applied to neural network topologies characterized by greater biological realism, and that long-range inhibitory connections can facilitate the generalization capacity of the process.

Keywords: artificial neural networks, self-organization, Hebbian learning, self-modeling, complex adaptive systems, Hopfield networks, artificial life, computational neuroscience

1. INTRODUCTION

The brain consists of a vast number of interacting elements. An important research question is how this complex adaptive system manages to give rise to large-scale coordination in the service of cognition, especially in the absence of a central controller or explicit knowledge of what would be the best neural connectivity. A promising approach is therefore the study of self-organization in artificial neural networks. Watson et al. (2011b) developed a self-optimization algorithm in Hopfield neural networks able to form associative memory of its attractor configurations through unsupervised learning of the Hebbian variety. This causes the networks to begin to reorganize their connectivity in a direction that minimizes the neural coordination constraints posed by the initial network architecture.

Previous work with this algorithm has been done using fully-connected networks, but without self-connections, and only with non-directed connections constrained to symmetric weights that are assigned in a random or highly modular manner (Watson et al., 2011a,c). More recently,

self-optimization has also been demonstrated in the case of continuous activation functions (Zarco and Froese, 2018a,b). This shows that the self-optimization process might be more generally applicable. Nevertheless, a concern with this work is that these network topologies are too artificial compared with those of actual neural networks. Accordingly, we propose that it would be more meaningful to employ the connectome of a real organism in order to better assess the scope of self-optimization.

A particularly suitable connectome comes from the nematode worm, *Caenorhabditis elegans*. This worm is one-millimeter-long and consists of only 959 cells, of which 302 belong to the nervous system. *C. elegans* is relevant in this research because it is a reference model in biology (White et al., 1986; Walker et al., 2000; Girard et al., 2006). It was the first multicellular organism whose genome has been sequenced in its entirety, as well as the first animal whose neural connections, called connectome, has been completed. *C. elegans* has also been studied in the field of artificial life using agent-based modeling (Izquierdo and Beer, 2015; Izquierdo, 2018).

In recent work, we demonstrated self-optimization in the *C. elegans* connectome (Morales and Froese, 2019), by turning it into a Hopfield neural network that captures the connectome's directed multigraph topology including its self-connections. We set two simulation experiments: (1) we ran the self-optimization algorithm with only excitatory (positive) connections, and (2) with 30% inhibitory (negative) connections arbitrarily assigned in a homogeneous fashion at the beginning of the algorithm. Under these conditions the *C. elegans* connectome showed a tendency to optimize its own connectivity, but more so in case (1). The addition of inhibitory synapses increased the difficulty of learning to find attractors with optimal neural coordination, and there remained a broader spread of attractors even after convergence. We hypothesize that this has to do with how coordination happens in functionally related neurons within clusters of the connectome.

Here we explore the possibility that this poor performance can be overcome by making inhibitory connections more concentrated between clusters, thereby also making our analysis more biologically plausible. We ran the self-optimization procedure in the whole *C. elegans* connectome, but also separately for each of the hierarchically organized functional clusters. We performed two sets of simulation experiments: (1) we arbitrarily assigned 30% inhibitory connections to local connections within each cluster, and ran self-optimization on each of the clusters as an independent network, and (2) we applied 30% of inhibitory connections to the whole connectome but restricted them to long-range inter-cluster connections, and ran the process on the entire connectome while also monitoring neural coordination within clusters.

The key finding of these simulation experiments is that the poor performance associated with the introduction of inhibitory connections can be successfully overcome by focusing inhibition to connections between clusters. This is the case both in terms of the number of attractors found and their energy levels: the process tends to converge on a more refined set of more optimal attractors, including attractors that normally would not be found by the network prior to self-optimization. Interestingly, while this

capacity to generalize to better attractors is also noticeable in the clusters when self-optimization is run on them independently, generalization is less marked when they are evaluated while embedded into the whole network—even though in the latter case they tend to converge on lower energy values because they do not have to overcome the added coordination constraints introduced by local inhibitory connections. This suggests that generalization to better attractor configurations is a property of the whole network, rather than being a simple reflection of generalization occurring at the level of local clusters.

2. METHODS

2.1. The Connectome

We ran the self-optimization algorithm in the connectome published by Jarrell et al. (2012). The database contains hermaphrodite neural system information (because males arise infrequently, at 0.1%), such as connection direction, type of connection (synapse or gap junction), and the number of connections between neurons. We translate the connectome into a directed multigraph, with neurons as nodes and connections as edges. Chemical synapses are modeled as single-directed links between neurons (for example, $A \rightarrow B$ indicates that neuron A is presynaptic to neuron B , and B is postsynaptic to A). Gap junctions are represented in the model as double-linked neurons (if two neurons, C and D , have a gap junction between them, there are two links: $C \rightarrow D$ and $D \rightarrow C$).

We assigned binary activation states $(-1, 1)$ to neurons. The number of connections between neurons was assigned as the weight of each edge, normalized in the interval $(0, 1)$. Both links in gap junctions were assigned the same weight, and values vary between 1 and 81 before normalization (and form a power law). Therefore, we clip to 1 the 15 high weight values, which we determine with an arbitrary cut-off of weights greater than 44. Reduction of this outliers broadens the state-space explorations during the self-optimization.

We did not also consider pharyngeal neurons because they belong to another independent neural system (Albertson and Thompson, 1976). Only 279 neurons are taken into account, with 5,588 connections. This differs from the number in our previous paper (282 neurons and 5,611 connections) because here we follow Sohn et al. (2011) in removing the neurons *VC6*, *CANR*, and *CANL* which do not have obvious connections.

Sohn et al. (2011) proposed a modular organization of the *C. elegans* connectome in five clusters based on a constraint community detection method for directed, weighted networks. This model shows hierarchical relationships between the clusters that define systemic cooperation between circuits with identified biological functions (mechanosensation, chemosensation, and navigation). This division also considers bilateral neural pairs present in the connectome so that the members of a pair should not be assigned to different structural clusters. There are two big clusters named 1 and 2. Smaller cluster names have hierarchical branch names: 1 (or 2) represents a big cluster branch in the left digit and small cluster branching is called 1 (or two rightward) in the right digit. **Table 1** shows the basic information of each cluster. The authors also observed many ties between the clusters

TABLE 1 | This table contains cluster information from the partition of Sohn et al. (2011), including the number of nodes and edges, average degree, and average shortest path of each cluster.

Cluster name	No. nodes	No. edges	Average node outgoing degree	Average shortest path	Cluster learning rate
Whole connectome	279	5,588 (3,392 intra, 2,196 inter-cluster)	20	2.5	0.00001
11	57	665	11.6	2.17	0.0000843
12	79	1,107	14	2.09	0.00005
13	14	115	8.2	1.52	0.0005
21	74	1,109	14.9	1.97	0.00005
22	55	396	7.2	3.08	0.0001416
11 + 12 + 13	150	2,704	18	2.23	0.0000207
21 + 22	129	1,980	15.3	2.34	0.0000283

Cluster names have hierarchical branch information: 1 (or 2) represents a former branch in the left digit and later branching is called 1 (or two rightward) in the right digit. First, we include information about the whole connectome before the partition, including the number of inter-cluster connections and intra-cluster ones. Then, we include information about the main 5 clusters. Finally, we include also information of the big clusters formed hierarchically from the five main clusters.

depended on hierarchical proximity. Cluster 11, 12, and 13 comprise a big cluster, and cluster 21 and 22 formed another grand cluster.

2.2. Model Dynamics

Asynchronous state updates are calculated with the following equation:

$$s_i(t+1) = \theta \left[\sum_j^N \left(\sum_k w_{ijk} \right) s_j(t) \right] \quad (1)$$

where s_i is the state of neuron i and w_{ijk} in the connection weight between neuron i and neuron j with index k (more than one tie with the same direction could arise between i and j). In a Hopfield network, a node i satisfies a constraint with its interaction with node j with index k if $s_i s_j w_{ijk} > 0$. System energy represents the constraint satisfaction level in the network:

$$E = - \sum_{ijk}^N w_{ijk}^O(t) s_i(t) s_j(t) \quad (2)$$

where w_{ijk}^O is the original weight configuration of w_{ijk} , the Hebbian learning changes during the process are managed in another variable.

The self-optimization algorithm consists on the repeating the following sequence of steps, each repetition is called a reset-convergence cycle:

1. Arbitrary assignment of states for the neurons (reset).
2. Convergence of the network for a certain time period, most frequently resulting in an attractor.
3. Application of Hebbian learning.

2.3. Introducing Inhibitory Connections

Morales and Froese (2019) explored two different weight configurations with self-optimization: when all connections are excitatory (positive), and when 30% are inhibitory (negative). In order to make the model more realistic, we introduced the inhibitory connections in the second weight configuration (Capano et al., 2015). This is because inhibitory connections are known to have an impact on network dynamics (Brunel, 2000). We found that the network shows a tendency to self-optimize when all connections are excitatory, but the 30% inhibitory connections restrict coordination and constraint satisfaction. Adding inhibitory connections will always have the effect of increasing the difficulty of constraint satisfaction, but it is also likely that this decrease in performance has to do with the fact that we distributed the inhibitory connections in a random fashion without taking the structural organization of the connectome into account. Therefore, we investigated the extent of self-optimization within each of the connectome's functional clusters with 30% inhibitory connections, and also self-optimization of the whole connectome when those inhibitory connections are concentrated between clusters.

More specifically, we run two sets of experiments: (1) self-optimization is run in each isolated cluster separately, and (2) we test for self-optimization in the whole connectome with inhibitory edges only assigned to inter-cluster connections and we monitor each embedded cluster. Since self-optimization in the network is sensitive to its size, we adjusted the learning rate in each isolated cluster in order to make the comparison fairer (see Table 1 for the learning rates). Python code of this simulation is available on GitHub¹.

3. RESULTS

Each experiment consists on the following setup (averaged from 10 different experiments with a different initial random number seed): the network is set to an initial configuration with only positive values and then we performed 1,000 reset-convergence cycles without Hebbian learning. Then, self-optimization is applied using 1,000 reset-convergence cycles that include Hebbian learning. Finally, another 1,000 reset-convergence cycles are applied without Hebbian learning using the learnt configuration obtained so far in order to show its stability. Note that these structural changes accumulated during learning are not directly reflected in the resulting figures. All the energy results shown in the figures were obtained by testing state configurations against the original connectome topology, because this reveals the extent to which the process was able to satisfy the original network constraints.

The experiment shown in Figure 1 explored self-optimization capacity in each isolated cluster, including the big clusters consisting of the join of smaller clusters. Each network tested separately show a tendency to self-coordinate during Hebbian learning, presenting a greater diversity of attractors. Some

¹https://github.com/aehecatl/self_opt_c_elegans

generalization capacity can also be seen, when a network starts to converge on energy values that were not previously seen during the initial phase. There are two exceptions: cluster 11 converges on a good energy value but one that was already included in the original distribution of energy values, and cluster 22 only converges on an average energy value of the ones previously encountered.

Figure 2 shows the experiments with 30% inhibitory connections arbitrarily assigned to only inter-cluster connections. We again find a tendency of the energy to decrease and the network to self-optimize, but the capacity for generalization to better previously unseen attractors is less notable. Nevertheless, the embedded clusters converge on better energy values compared to the isolated clusters, although

this may be partially because the inhibitory connections were moved to the inter-cluster domain, thereby also decreasing the difficulty of intra-cluster coordination. However, we know that this decrease in intra-cluster complexity is not the whole story because there is one exception: cluster 13 performs worse under these embedded conditions compared to isolated conditions.

This leads us to ask about the performance of self-optimization at the level of the whole connectome. **Figure 3** shows that restricting inhibitory connections to the inter-cluster domain has the effect of facilitating the self-optimization process: it now consistently generalizes to a more refined set of energy values that are much lower. This occurs despite the fact that both conditions feature the same overall number of inhibitory connections.

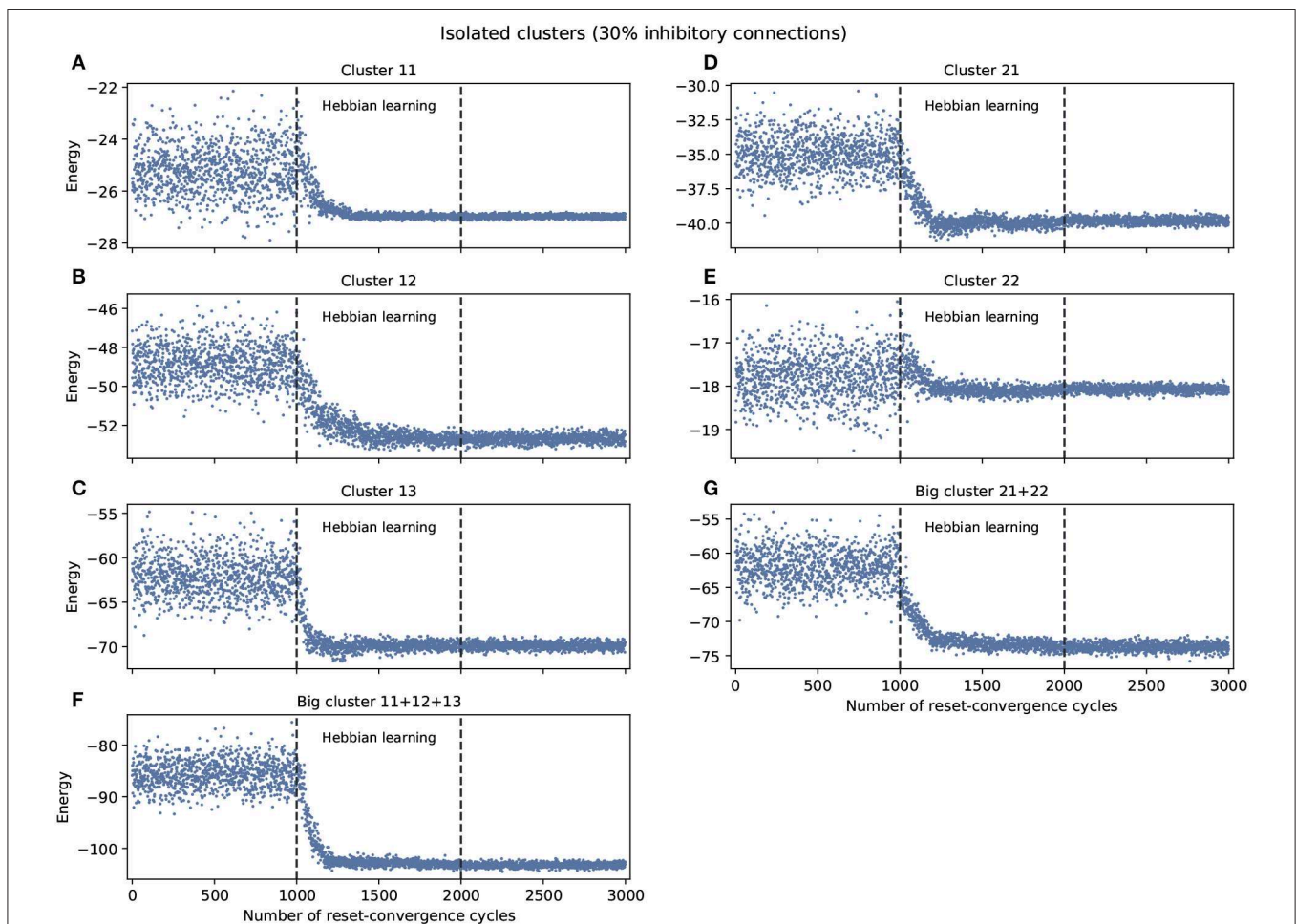
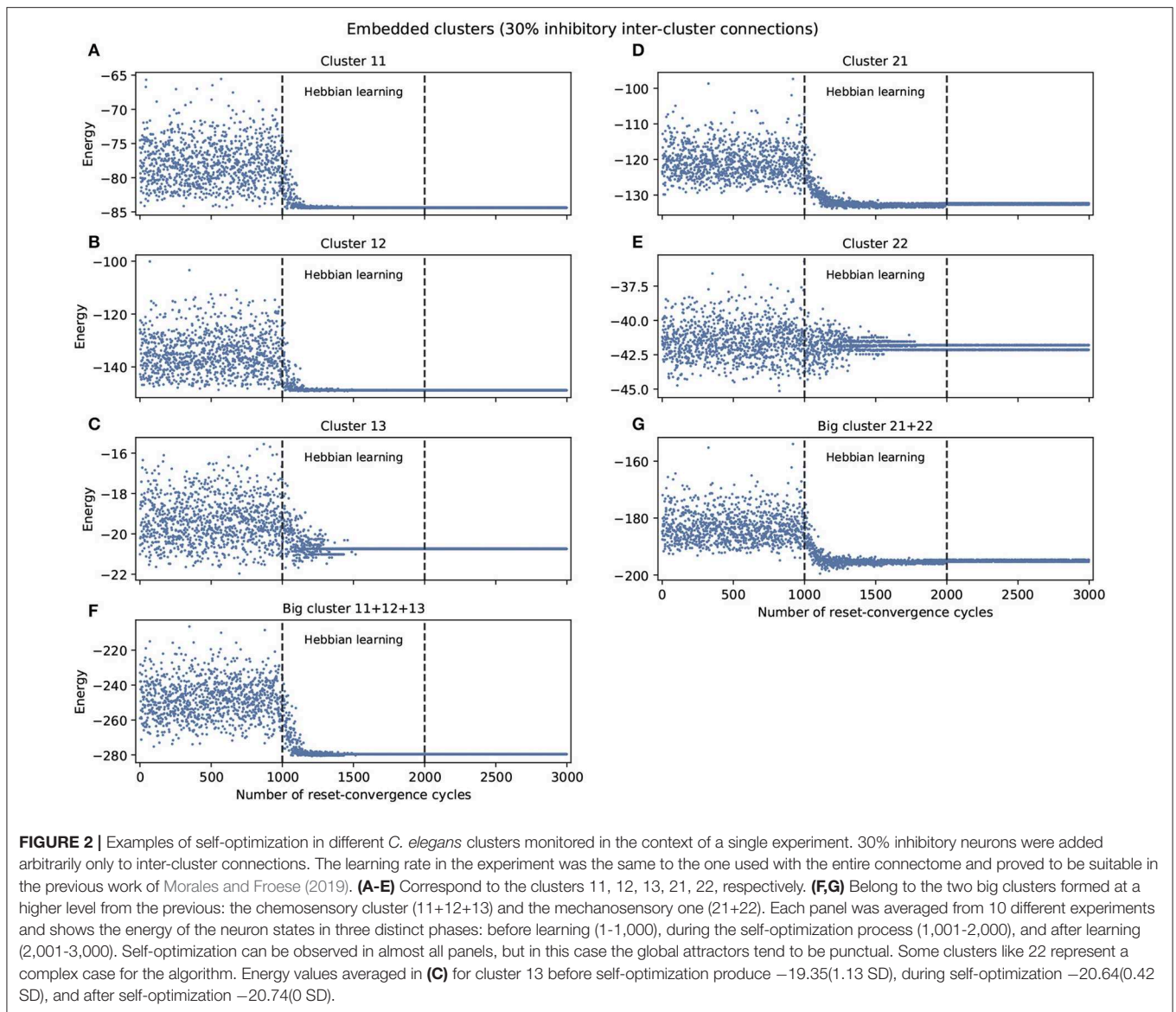


FIGURE 1 | Examples of self-optimization in different *C. elegans* clusters with 30% inhibitory connections; each panel was run separately (independent to the rest of the connectome). The learning rate in each experiment was proportional (regarding the edges) to the one used with the entire connectome and proved to be suitable in the previous work of Morales and Froese (2019). (A–E) Correspond to the clusters 11, 12, 13, 21, 22, respectively. (F,G) Belong to the two big clusters formed at a higher level from the previous: the chemosensory cluster (11 + 12 + 13) and the mechanosensory one (21 + 22). Each panel was averaged from 10 different experiments and shows the energy of the neuron states in three distinct phases: before learning (1–1,000), during the self-optimization process (1,001–2,000), and after learning (2,001–3,000). Self-optimization can be observed in almost all panels, but tend to remain a diversity of attractors. The difference in y-scale of each panel underline the complexity of the problem to be solved by self-optimization. Energy values averaged in (A) before self-optimization produce -25.22 (0.91 SD), during self-optimization -26.71 (0.61 SD), and after self-optimization -26.98 (0.06 SD). In the case of (E) we have -17.78 (0.5 SD), -18.01 (0.24 SD), and -18.07 (0.07 SD), respectively.

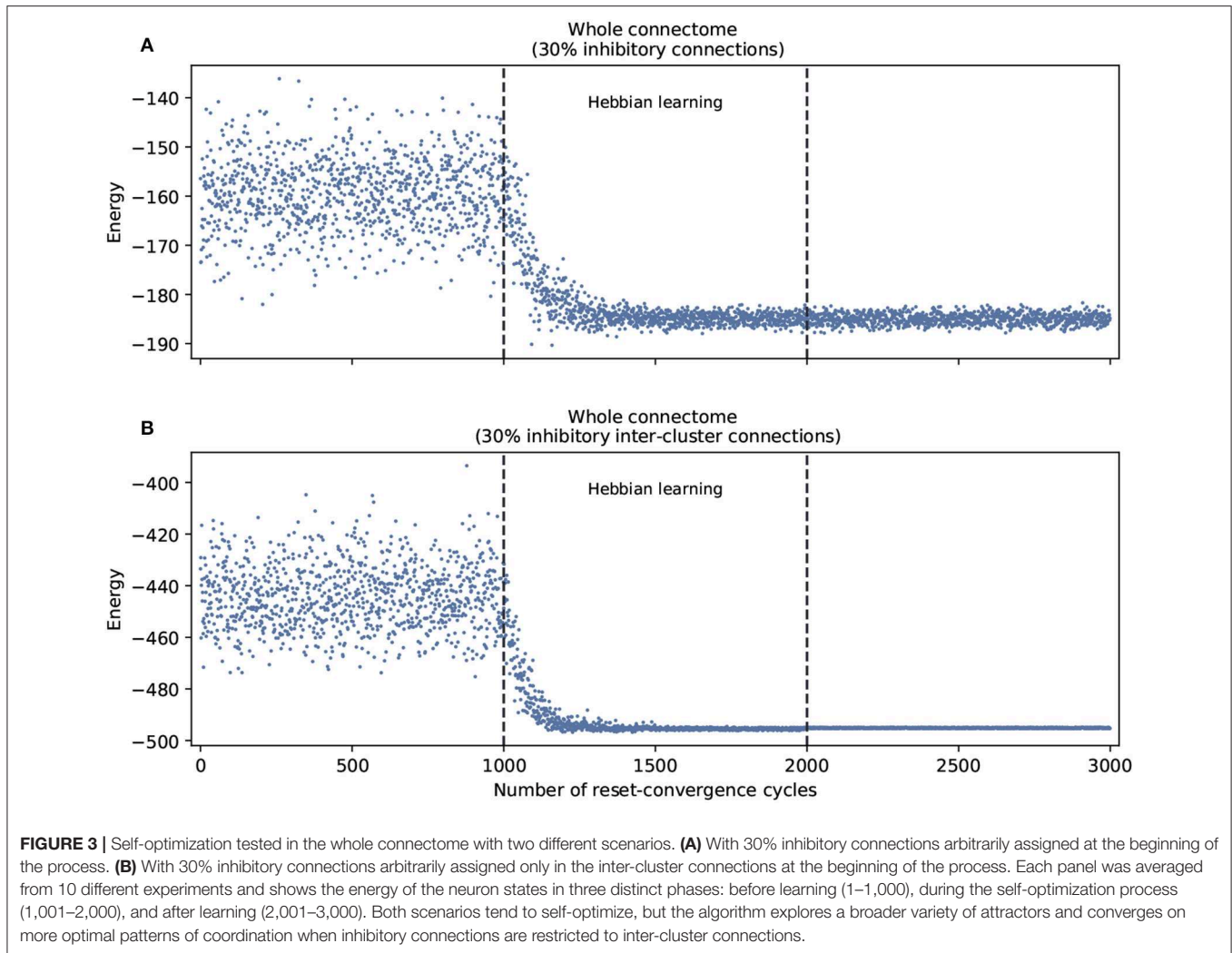


4. DISCUSSION

We successfully demonstrated the capacity of self-optimization for the case of the *C. elegans* connectome. Through repeated reset-convergence cycles, the network managed to generalize to previously unseen attractors with better coordination constraint satisfaction. Moreover, we managed to improve on previous work by showing that inhibitory connections do not hinder this process as long as they are concentrated to connections between clusters.

For simplicity, we assigned all inhibitory connections to inter-cluster connections in an arbitrary way. However, in real neural networks it is whole neurons, not isolated connections, that are inhibitory. Future work could therefore further improve the biological realism of our model by taking into account the excitatory or inhibitory functions of the neurotransmitters associated with each of the neurons in the connectome (Riddle et al., 1997; Pereira et al., 2015).

We also note that here we only explored the dynamics of the network in an uncoupled mode. Accordingly, an outstanding challenge is to embed the model of the connectome in whole worm simulations to explore the relationship between coupled and uncoupled dynamics (Izquierdo and Bührmann, 2008; Zarco and Froese, 2018b). So far it is unknown whether self-optimization can also occur when the network is in a coupled mode. Nevertheless, it has been speculated that the uncoupled mode of self-optimization could reflect the prevalent need for sleep among animals (Woodward et al., 2015). If this is on the right track, our model could be developed into a scientific hypothesis to inform current debates about the function of the quiescent state observed in *C. elegans* (Raizen et al., 2008; Trojanowski and Raizen, 2016). Future modeling work could also explore similarities and differences between this proposal and other neural network models of the function of sleep (Hopfield et al., 1983; Fachechi et al., 2019).



One limitation of our work is that the model is not sufficiently realistic compared with living systems and their complex interactions at different levels. We can overcome this limitations by implementing our model under different attractor dynamics like heteroclinic or slow and fast dynamics in synapses (Izhikevich, 2007).

DATA AVAILABILITY STATEMENT

The raw data supporting the conclusions of this article will be made available by the authors, without undue reservation, to any qualified researcher.

REFERENCES

- Albertson, D. G., and Thompson, J. (1976). The pharynx of *Caenorhabditis elegans*. *Philos. Trans. R. Soc. Lond. Biol. Sci.* 275, 299–325. doi: 10.1098/rstb.1976.0085
- Brunel, N. (2000). Dynamics of sparsely connected networks of excitatory and inhibitory spiking neurons. *J. Comput. Neurosci.* 8, 183–208. doi: 10.1023/A:1008925309027

AUTHOR CONTRIBUTIONS

All authors listed have made a substantial, direct and intellectual contribution to the work, and approved it for publication.

ACKNOWLEDGMENTS

We thank OIST Scientific Computing & Data Analysis Section for the use and support of the High Performance Computing Cluster. We thank Thomas Burns for helpful discussion.

- Capano, V., Herrmann, H. J., and de Arcangelis, L. (2015). Optimal percentage of inhibitory synapses in multi-task learning. *Sci. Rep.* 5:9895. doi: 10.1038/srep09895
- Fachechi, A., Agliari, E., and Barra, A. (2019). Dreaming neural networks: forgetting spurious memories and reinforcing pure ones. *Neural Netw.* 112, 24–40. doi: 10.1016/j.neunet.2019.01.006

- Girard, L. R., Fiedler, T. J., Harris, T. W., Carvalho, F., Antoshechkin, I., Han, M., et al. (2006). Wormbook: the online review of *Caenorhabditis elegans* biology. *Nucl. Acids Res.* 35, D472–D475. doi: 10.1093/nar/gkl894
- Hopfield, J. J., Feinstein, D., and Palmer, R. (1983). ‘unlearning’ has a stabilizing effect in collective memories. *Nature* 304, 158–159. doi: 10.1038/304158a0
- Izhikevich, E. M. (2007). *Dynamical Systems in Neuroscience*. Cambridge, MA: MIT Press. doi: 10.7551/mitpress/2526.001.0001
- Izquierdo, E., and Beer, R. (2015). An “integrated neuromechanical model of steering in *C. elegans*,” in *The Thirteenth European Conference on Artificial Life*, eds P. Andrews, L. Caves, R. Doursat, S. Hickinbotham, F. Polack, S. Stepney, T. Taylor, and J. Timmis (Cambridge, MA: MIT Press), 199–206.
- Izquierdo, E., and Bührmann, T. (2008). “Analysis of a dynamical recurrent neural network evolved for two qualitatively different tasks: walking and chemotaxis,” in *The eleventh international conference on the simulation and synthesis of living systems*, eds S. Bullock, J. Noble, R. Watson, and M. Bedau (Cambridge, MA: MIT Press), 257–264.
- Izquierdo, E. J. (2018). Role of simulation models in understanding the generation of behavior in *C. elegans*. *Curr. Opin. Syst. Biol.* 13, 93–101. doi: 10.1016/j.coisb.2018.11.003
- Jarrell, T. A., Wang, Y., Bloniarz, A. E., Brittin, C. A., Xu, M., Thomson, J. N., et al. (2012). The connectome of a decision-making neural network. *Science* 337, 437–444. doi: 10.1126/science.1221762
- Morales, A., and Froese, T. (2019). “Self-optimization in a hopfield neural network based on the *C. elegans* connectome,” in *The 2019 Conference on Artificial Life*, eds H. Fellerman, J. Bacardit, A. Goñi-Moreno, and R. Fuchsln (Cambridge, MA: MIT Press), 448–453.
- Pereira, L., Kratsios, P., Serrano-Saiz, E., Sheftel, H., Mayo, A. E., Hall, D. H., et al. (2015). A cellular and regulatory map of the cholinergic nervous system of *C. elegans*. *Elife* 4:e12432. doi: 10.7554/eLife.12432
- Raizen, D. M., Zimmerman, J. E., Maycock, M. H., Ta, U. D., You, Y.-J., Sundaram, M. V., et al. (2008). Lethargus is a *Caenorhabditis elegans* sleep-like state. *Nature* 451:569. doi: 10.1038/nature06535
- Riddle, D. L., Blumenthal, T., Meyer, B. J., and Priess, J. R., (eds.). (1997). *C. elegans II. 2nd Edn*. Cold Spring Harbor, NY: Cold Spring Harbor Laboratory Press.
- Sohn, Y., Choi, M.-K., Ahn, Y.-Y., Lee, J., and Jeong, J. (2011). Topological cluster analysis reveals the systemic organization of the *Caenorhabditis elegans* connectome. *PLoS Comput. Biol.* 7:e1001139. doi: 10.1371/journal.pcbi.1001139
- Trojanowski, N. F., and Raizen, D. M. (2016). Call it worm sleep. *Trends Neurosci.* 39, 54–62. doi: 10.1016/j.tins.2015.12.005
- Walker, D. W., McColl, G., Jenkins, N. L., Harris, J., and Lithgow, G. J. (2000). Natural selection: evolution of lifespan in *C. elegans*. *Nature* 405:296. doi: 10.1038/35012693
- Watson, R. A., Buckley, C. L., and Mills, R. (2011a). Optimization in “self-modeling” complex adaptive systems. *Complexity* 16, 17–26. doi: 10.1002/cplx.20346
- Watson, R. A., Mills, R., and Buckley, C. L. (2011b). Global adaptation in networks of selfish components: emergent associative memory at the system scale. *Artif. Life* 17, 147–166. doi: 10.1162/artl_a_00029
- Watson, R. A., Mills, R., and Buckley, C. L. (2011c). Transformations in the scale of behavior and the global optimization of constraints in adaptive networks. *Adapt. Behav.* 19, 227–249. doi: 10.1177/1059712311412797
- White, J. G., Southgate, E., Thomson, J. N., and Brenner, S. (1986). The structure of the nervous system of the nematode *Caenorhabditis elegans*. *Philos. Trans. R. Soc. Biol. Sci.* 314, 1–340. doi: 10.1098/rstb.1986.0056
- Woodward, A., Froese, T., and Ikegami, T. (2015). Neural coordination can be enhanced by occasional interruption of normal firing patterns: a self-optimizing spiking neural network model. *Neural Netw.* 62, 39–46. doi: 10.1016/j.neunet.2014.08.011
- Zarco, M., and Froese, T. (2018a). Self-modeling in Hopfield neural networks with continuous activation function. *Procedia Computer Science* 123, 573–578. doi: 10.1016/j.procs.2018.01.087
- Zarco, M., and Froese, T. (2018b). Self-optimization in continuous-time recurrent neural networks. *Front. Robot. AI* 5:96. doi: 10.3389/frobt.2018.00096

Conflict of Interest: The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2020 Morales and Froese. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.



New Methods for the Steady-State Analysis of Complex Agent-Based Models

Chico Q. Camargo*

Oxford Internet Institute, University of Oxford, Oxford, United Kingdom

Among all tools used to understand collective human behavior, few tools have been as successful as agent-based models (ABMs). These models have been particularly effective at describing emergent social behavior, such as spatial segregation in neighborhoods or opinion polarization on social networks. ABMs are particularly common in the study of opinion and belief dynamics, being used by fields ranging from anthropology to statistical physics. These models, much like the social systems they describe, often do not have unique output variables, scales, or clear order parameters. This lack of clearly measurable emergent behavior makes such complex ABMs difficult to study, ultimately limiting their application to cases of empirical interest. In this paper, we introduce a series of approaches to analyze complex multidimensional ABMs, drawing from information theory and cluster analysis. We use these approaches to explore a multi-level agent-based model of ideological alignment introduced by Banisch and Olbrisch to extend Mäs and Flache's argument communication theory of bi-polarization. We use the tools introduced here to perform a thorough analysis of the model for small system sizes, identifying the convergence toward steady-state behavior, and describing the full spectrum of steady-state distributions produced by this model. Finally, we show how the approach we introduced can be easily adapted for larger implementations, as well as for other complex agent-based models of social behavior.

Keywords: complex systems, agent-based modeling, computational social science, opinion dynamics, belief dynamics, social influence, polarization, cognitive-evaluative maps

OPEN ACCESS

Edited by:

Carlos Gershenson,
National Autonomous University of
Mexico, Mexico

Reviewed by:

Sven Banisch,
Max Planck Institute for Mathematics
in the Sciences, Germany
Feng Fu,
Dartmouth College, United States

*Correspondence:

Chico Q. Camargo
chico.camargo@oii.ox.ac.uk

Specialty section:

This article was submitted to
Interdisciplinary Physics,
a section of the journal
Frontiers in Physics

Received: 31 October 2019

Accepted: 18 March 2020

Published: 08 April 2020

Citation:

Camargo CQ (2020) New Methods for
the Steady-State Analysis of Complex
Agent-Based Models.
Front. Phys. 8:103.
doi: 10.3389/fphy.2020.00103

1. INTRODUCTION

Over the last decades, computational social science has risen as a strongly empirical discipline, drawing on data science methods to tackle high-dimensional large data sets that cannot be understood with simple analytical tools. This is particularly true in the study of public attention and public opinion dynamics: there are multiple studies looking at large-scale trends in Internet search queries, online petitions, and forums, as well as applying natural language processing methods to news articles and social media activity. It is now possible to quantify, to a degree, what people care about, and what they think of it.

This increasing availability of data on individual and public opinion calls for realistic, theory-informed models. Models of opinion change and belief dynamics have traditionally been studied by a large number of disciplines, including but not limited to economics [1–3], political science [4], sociology [5], anthropology [6], philosophy [7], and psychology [8] among others. There is also a tradition in statistical physics, dating back to the voter model [9–12], but also considering models such as the majority rule model [13], the Sznajd model [14], and a number of bounded confidence models [15]. Each model typically describes opinion change through

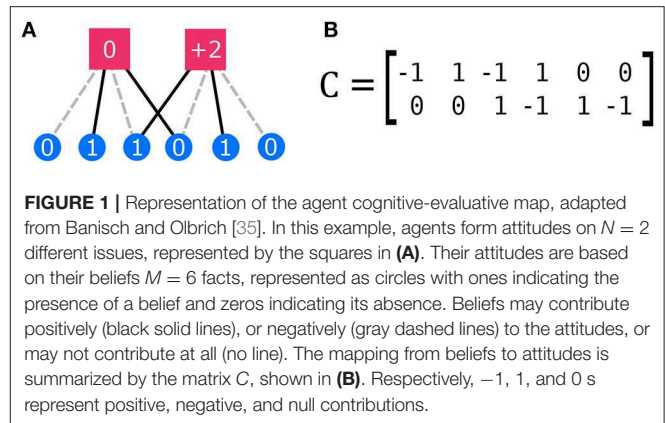
a fixed strategy, where an agent might update their beliefs to more accurate values [16–18], or perhaps might seek conformity by following either the majority around them [13], or by copying the mean opinion [4]. The effects of social influence might vary with the distance between one’s own opinion and the advocated one [19–21], on the details of how the new information is presented [22] or even to meta-information [23–25]. Opinion dynamics models often also take into account the structure of the social networks where agents are embedded. In these models, opinion formation is often described as a result of the combination of social structure and behavior, as agents in different parts of a social network will be exposed to different sources of information, while the social network itself might change over time, as agents choose to change their own connections according to the behavior and opinions of their neighbors [11, 26, 27].

Rather than producing an exhaustive list of models and modeling choices, this study aims to develop methods that allow for a thorough exploration of complex models. Many of the models presented above, much like the social systems they describe, have multiple output variables, often displaying divergent behavior in one coarse-graining scale while displaying convergent behavior in another. This makes such complex models hard to study, particularly as their application into questions of empirical interest requires expanding models to multidimensional landscapes and parameter spaces, where emergence and convergence are difficult to identify in first place.

With these problems in mind, in this paper we introduce a series of tools that provide a scalable way to explore the parameter space of complex agent-based models. As a case study, we use an multi-level opinion dynamics agent-based model which contains all of these features—no clear output variable, multidimensional parameter space and output space. We perform a thorough analysis of the model for small system sizes, and show how the same analysis could be performed for larger implementations of other complex models.

The agent-based model we use as a case study was originally introduced as a toy model of opinion polarization. In recent years, the spread of information on social networks has been described a driving force behind political polarization, through mechanisms of homophily leading to “filter bubbles” or “echo chambers” [11, 28, 29]. While a more robust body of evidence is needed to clarify the many roles of online social networks in political polarization, the role of homophily and social influence in the process of opinion polarization is already well-described by concepts such as the argument communication theory of bi-polarization, introduced by Mäs and Flache [30]. This theory proposes to account for the emergence of a bi-modal distribution in opinions through an amplification of small differences between individuals. It draws from the theory of informational influence, or persuasive argument theory [16–18], while assuming that homophily with respect to an individual’s opinions [28, 31–34] will be the main factor behind communication and opinion change. As argued by Mäs and Flache, the cognitive-social bias of homophily is a sufficient mechanism to account for the emergence of a bimodal opinion distribution.

The simplicity of Mäs and Flache’s theoretical model is also its limitation, in that it focuses on the emergence of polarization around a single issue, or a single pair of opposing issues on



an axis. This limitation has been addressed by an extended computational model proposed by Banisch and Olbrich [35], who draw from structural theories of attitude dynamics [36–38] to make a distinction between individual beliefs held by an agent and their attitudes toward multiple issues. The relations between beliefs and attitudes are framed by Banisch and Olbrich as cognitive-evaluative maps shared by a population [39]. In their computational model, an individual’s beliefs are encoded as a vector x of binary values, while their attitudes are represented as another vector y , this one with integer values, which depend on the belief vector but also on a cognitive-evaluative matrix C , through the linear equation $y = C \cdot x$ (in the notation used here). In the example shown in **Figure 1A**, a total of six beliefs determines an agent’s attitude toward two issues. Each issue is affected positively by two issues, negatively by two others, and is not affected by the last two. This can be represented as a bipartite graph where every belief is connected to an opinion, which can be summarized by the adjacency matrix shown in **Figure 1B**. The separation between belief and attitude makes this network different from other network models of belief dynamics [40], where beliefs affect each influence each other directly. In this two-level model, in principle, unless two agents interact, one agent does not have access to another agent’s beliefs—while the decision to interact might be based on attitudes only.

As described above, Banisch and Olbrich’s model of ideological alignment is a multi-dimensional agent-based model which can exhibit emergent behavior in more than one level, posing an interesting challenge for current methods of analysis of ABMs. In the next sections, we explore the parameter space of this model by investigating the ensemble of all cognitive-evaluative maps for systems with small numbers of beliefs and attitudes. We then introduce different approaches to analyze the convergence of the model, as well as the range of steady-states it can produce. Finally, we argue how these approaches can easily be applied to other complex agent-based models of social behavior.

2. METHODS

2.1. Multi-Level Agent Based Model

Following Banisch and Olbrich’s model [35], we define the cognitive-evaluation matrix relating M beliefs to N issues as a $N \times M$ matrix C . We limit entries c_{ij} to values within

$\{-1, 1, 0\}$, corresponding to the attitude toward an issue j being affected negatively, positively, or unaffected by a belief i . This implies a total of 3^{MN} possible cognitive-evaluative matrices. The exponential growth with M and N is a product of the combinatorial nature of the problem, since *a priori* the relation between a pair of beliefs i_1 and j_1 does not impose any constraint on the relation between any pair i_2 and j_2 . Consequently, the number of possible C matrices quickly grows beyond what would be effectively enumerable. For $M = 2$, $N = 2$, there is a total of $3^4 = 81$ possible matrices, while for larger systems such as $M = 10$ beliefs affecting $N = 3$ issues, this number grows to $3^{30} \approx 2 \times 10^{14}$. When considering the output of every agent based model, we take into account how multiple matrices might be equivalent under symmetry operations. These operations, which represent all permutations of an agent's beliefs and opinions (e.g., replacing belief i for belief j), result in a smaller set of isomorphic graphs connecting beliefs to opinions, thus mitigating the exponential growth described above. In this brief study, we focus on three case studies where a thorough study of the matrix ensemble is possible, once such symmetries are taken into consideration: namely $M = 4$, $N = 2$ and $M = 3$, $N = 5$.

For every matrix C in each matrix ensemble, we run a total of 20 simulations with different random seeds. For every random seed, we initialize 1,000 agents with random sets of beliefs, i.e., initializing their beliefs as random vectors $x \in \{0, 1\}^M$, and mapping them to $y \in \mathbb{Z}^N$ attitude vectors through $y = C \cdot x$. We then iterate every simulation through 15,000 time steps, which we show is enough for model convergence. In every time step, for every agent a_1 in the model, we select another agent a_2 at random, measure the homophily between them, and if this homophily is above a given threshold, agent a_1 selects a random belief from a_2 's beliefs, and copies it. With 1,000 agents and 15,000 time steps, every simulation runs for a total of 15×10^6 interactions between agents.

As described above, the similarity or homophily between agents can be measured in multiple ways. In this study, we define homophily as measured by one minus the normalized Manhattan distance between the attitudes of a pair of agents. The reasons for this choice are many. First, if the distance between agents were to be measured in belief space, i.e., according to their belief vectors only, the dynamics of the agent based model would be trivial: agents would move toward each other, aggregating in a few points in belief space, and no other kind of dynamics would be possible. In other words, agents would concentrate in a finite number of $x \in \{0, 1\}^M$ points in belief space.

This kind of dynamics corresponds to a series of bounded confidence models in opinion dynamics, such as the ones introduced by Deffuant et al. [15] and by Krause and Stöckler [41] and Hegselmann and Krause [42], both of which were expanded by many following works [43–48]. In this category of models, for high enough homophily thresholds, agents might cluster in a few points, while for lower thresholds they would eventually all collapse into a single set of beliefs x , depending on the initial distribution of a agents in the opinion space, but not depending on the cognitive-evaluative matrix C . If one were to assume, for example, that every set of beliefs is equally likely *a priori*, thus defining the initial conditions of the simulation as a uniform

distribution over $\{0, 1\}^M$, then every point in belief space would be equally likely to become the steady state of the system, upon a random perturbation to the uniform initial condition of the model. The corresponding dynamics in attitude space $y \in \mathbb{Z}^N$ would also be one of aggregation toward a few attitude vectors y , and the likelihood of convergence toward a specific attitude vector would be proportional to the fraction of x vectors that map to y through $y = C \cdot x$. Other than that, unless the combination of a particular initial distribution over belief space and the right homophily threshold could allow to the formation of two clusters, measuring homophily as a function of belief homophily would only lead to the formation of homogeneous steady states where all agent have exactly the same opinions and beliefs.

Given that the dynamics induced by any distance metric in belief space will inevitably lead to agents aggregation in both belief and attitude space, what is left is to investigate the dynamics produced by metrics in attitude space. For this choice, if one were to use the Euclidean or L2 norm to measure the distance between attitude vectors y , for a given set of $y_1 = (0, 0)$, $y_2 = (1, 1)$, and $y_3 = (0, 2)$, one would obtain $dist(y_1, y_2) < dist(y_1, y_3)$. Were one to use the Manhattan or L1 norm, they would find $dist(y_1, y_2) = dist(y_1, y_3)$. In the absence of reasons to argue that y_1 differs more from y_3 than from y_2 , we will pick the simplest assumption, and use the Manhattan norm for simplicity.

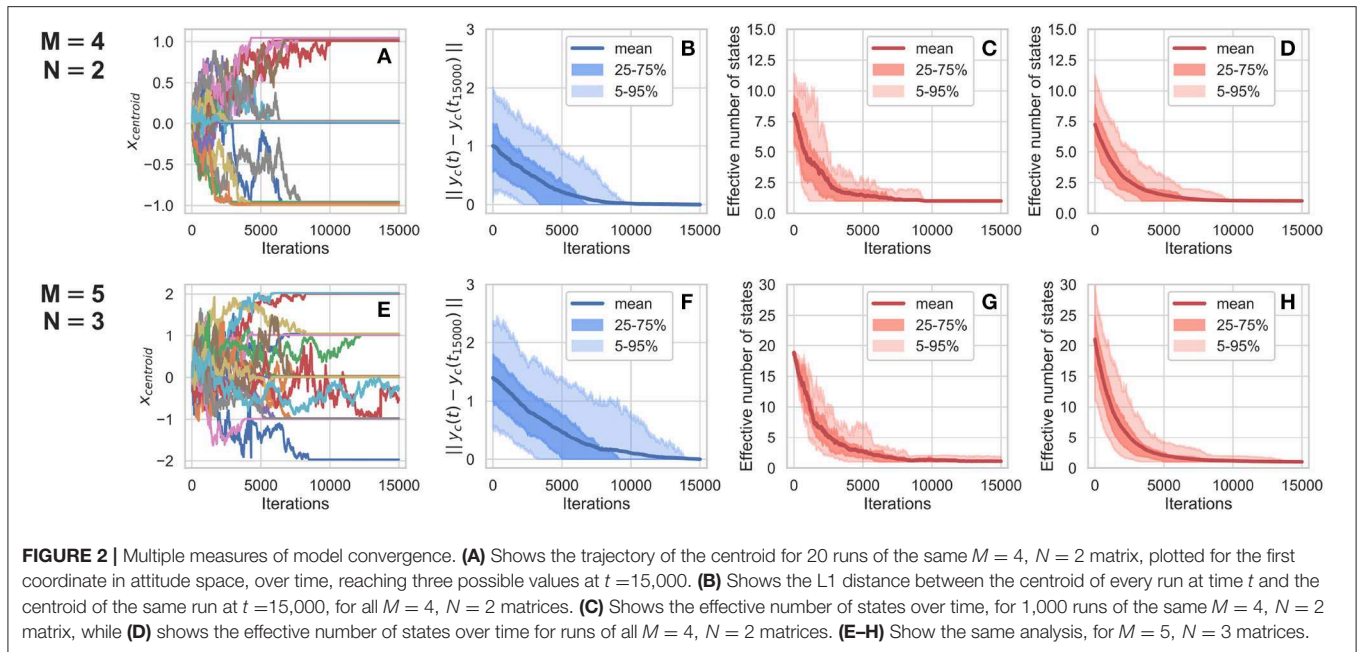
2.2. Measuring Simulation Outputs

To test whether simulations with the same C matrix but initialized with different random seeds might produce different steady-state distributions, we first assess the variability within multiple runs of the same matrix, as shown in **Figure 2**. Since every run of the model produces 1,000 trajectories over a M -dimensional belief space and a N -dimensional opinion space, we represent the state of an individual run over time with a set of summary statistics: its centroid $y_{\text{centroid}} \in \mathbb{R}^N$, its covariance matrix Σ_{ij} , and its maximum width in each of its principal axes, which can be identified by decomposing Σ_{ij} into its scaling and rotational components.

We measure the spread of agents over time for every run in two ways. First, we calculate the mean distance from the centroid of a simulation run at a given time step and the centroid of its steady state (i.e., its centroid after 15,000 steps). Second, we measure the effective number of states over time for each run. The effective number of states is a measure inspired in entropy-based measures of diversity, which have their origin in information theory [49, 50]. We define it as 2 to the power of the entropy of the distribution of agents over multiple states, as shown in Equation (1).

$$\text{ENS}(t) = 2^{-\sum_i p_i(t) \log_2 p_i(t)} \quad (1)$$

In Equation (1) above, the entropy term is summed over the proportion $p_i(t)$ of agents occupying state y_i at time t , for all states y_i in attitude space. In essence, the effective number of states is a measure of the diversity of sets of attitudes taken by the agents in a run at a given point in time: in other words, of how their attitudes are divided between multiple simultaneous states (or sets of attitude values y), weighed according to how



many agents adopt them at that point in time. This measure is highest when agents are evenly distributed across many states and lowest when they concentrated at a single state. In the social sciences, equivalent approaches have been used to describe the effective number of parties in a parliament [51], as well as the effective number of issues from a political agenda receiving public attention at the same time [52].

Finally, in addition to analyzing every simulation for every C matrix, we also cluster groups of steady-state distributions according to which points in attitude space are occupied at $t = 15,000$ by a given run.

3. RESULTS

In this section, we present the results of simulations for the three case studies mentioned above: $M = 4$, $N = 2$ and $M = 3$, $N = 5$. Unless specified, we use a distance threshold of $\beta = 1$. Unlike Banisch and Olbrich's study [35], which shows the results for a selected set of matrices that could produce a varied set of behaviors, we focus on the behavior emerging from the whole ensemble of matrices defined by a given (M, N) pair.

3.1. Studying Model Convergence

Figure 2 displays an analysis of convergence for this agent-based model, for system sizes of $M = 4$, $N = 2$, and $M = 3$, $N = 5$. In summary, it shows that different runs of the same matrix can produce different results, that convergence typically happens before $t = 15,000$ steps, and that this convergence is usually to a single point in attitude space.

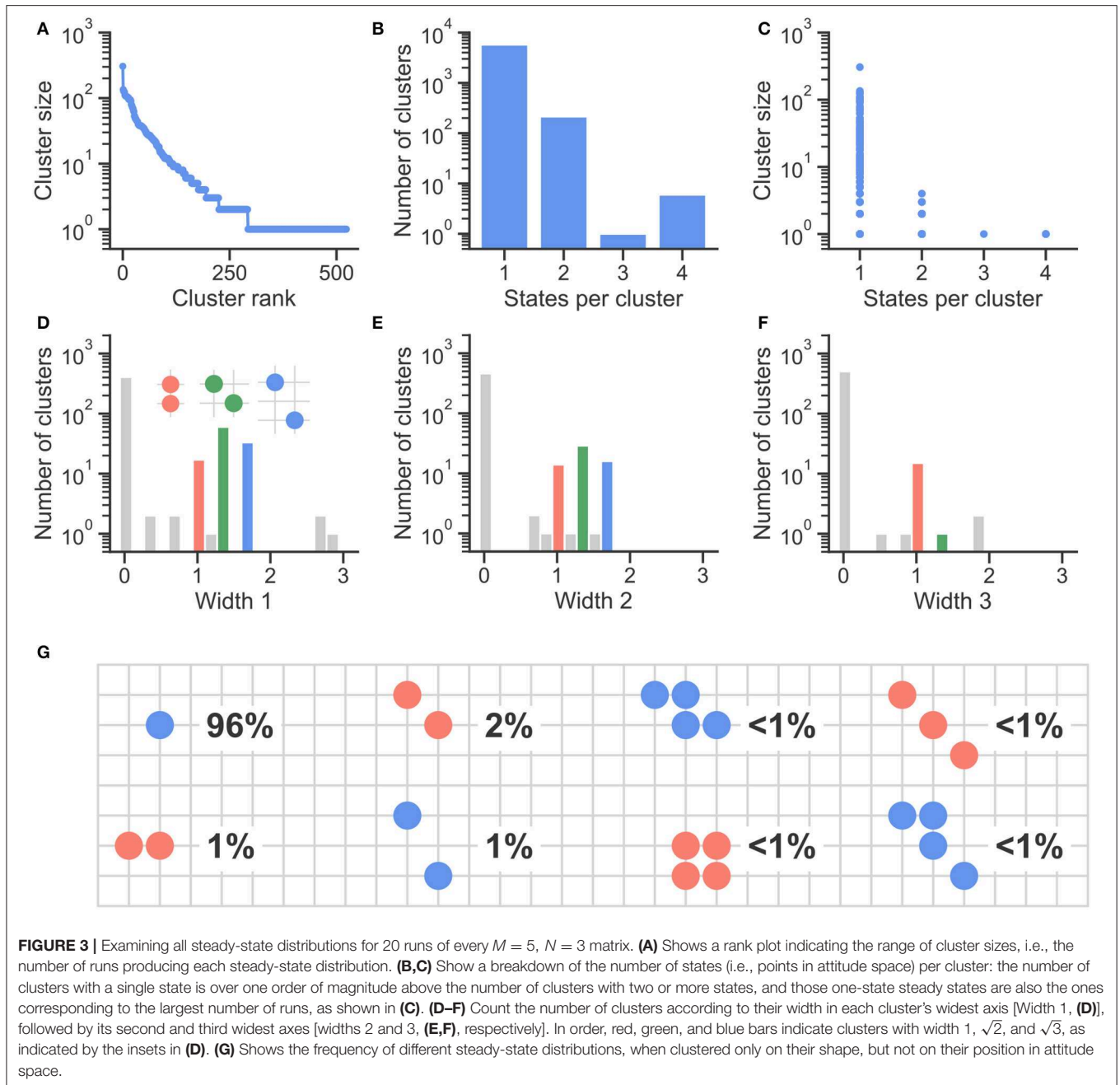
The figure compares the full ensembles of $M = 4$, $N = 2$ and $M = 5$, $N = 3$ matrices with the matrices $C_{2 \times 4}$ and $C_{3 \times 5}$

specified in Equation (2):

$$C_{2 \times 4} = \begin{bmatrix} -1 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{bmatrix} \quad C_{3 \times 5} = \begin{bmatrix} 0 & -1 & 1 & 1 \\ 1 & -1 & 0 & 1 \\ 0 & 0 & -1 & 0 \end{bmatrix} \quad (2)$$

The first point is illustrated by **Figures 2A,E**. Both panels show the trajectory of the first coordinate of the centroid of 20 different runs of the same C matrix. In plotting these time series, a small increment of 0.01 was added to the y -value of each run, to make visible the many horizontal lines that otherwise would be overlaid. The panels indicate that most centroids converge to a value before 15,000 steps, but that the value itself varies across runs. For this particular choice of C matrices, centroids stabilized at values of -1 , 0 , and $+1$ for $C_{2 \times 4}$, and -2 , -1 , 0 , $+1$, and $+2$ for $C_{3 \times 5}$. The fact that the values of 2 and -2 are not observed for $C_{2 \times 4}$ and that neither 3 or -3 is observed for $C_{3 \times 5}$ is likely due to these specific maps. Still, the diversity of centroid values presented in both panels is enough to show the kind of behavior that would be erased if one were to average multiple runs for the same C matrix.

Naturally, displaying the results of a single pair of matrices is no argument for general convergence. The model convergence around 15,000 time steps for these (M, N) pairs is further presented in all other panels in **Figure 2**. **Figures 2B,F** show the L1 distance between the centroid of the distribution of agents in attitude space at time t and the same distribution at time $t = 15,000$, averaged over all $M = 4$, $N = 2$ and $M = 5$, $N = 3$ matrices, respectively for each panel. Shaded areas represent the 25–75 and 5–95% intervals of the distribution of the distance to steady-state centroids, showing that the convergence observed at $t = 15,000$ is not an average phenomenon, and also not unique to $C_{2 \times 4}$ and $C_{3 \times 5}$, but rather that convergence is observed for both whole matrix ensembles.



The remaining panels show the evolution of the effective number of states over time, for 1,000 runs of $C_{2 \times 4}$ and $C_{3 \times 5}$ (**Figures 2C,G**) and for single runs of all matrices in that (M, N) pair (**Figures 2D,H**). The effective number of states, described in Equation (1), measures the diversity of points in attitude space occupied by the multiple agents in a model over time. In all panels, this effective number quickly converges to approximately 1.0, both on average and as a whole, as indicated by the shaded areas. This convergence implies that most runs ultimately lead to steady states occupying a single point in attitude space, for both $M = 4$, $N = 2$ and $M = 5$, $N = 3$ matrices.

3.2. Analyzing Steady States

In the previous section, we established that the model usually converges before 15,000 steps, that a typical run converges to a single point in attitude space, but that different runs of the same matrix might result in path-dependent symmetry breaking. In this section, we examine the range of steady-state distributions produced by multiple runs of this model for many C matrices, clustered according to which points in attitude space (i.e., which states) are occupied at $t = 15,000$ by a given run.

The results of the clustered by steady-state distributions are shown in **Figure 3** for 20 runs of every $M = 5$, $N = 3$ matrix.

Figure 3A shows a rank plot indicating the range of cluster sizes, i.e., the number of runs producing each steady-state distribution. As evidenced by the log-scale on the y axis, this is a long-tail distribution: most runs produce the same few steady-state distributions, while most steady states are only observed for 10 runs or less.

From **Figure 3B**, we see that most steady-state clusters correspond to single states, while the number of clusters with two or more states is over an order of magnitude smaller. **Figure 3C** compares the number of states with the number of runs falling into each cluster, i.e., the cluster size: it shows that most large clusters are single-state clusters, followed by two-state clusters.

As indicated by the top three panels, most runs of this model result in a few single-state clusters, while wider steady-state distributions correspond to a small proportion of all resulting steady-states of the model, with many distributions corresponding to only a few model runs each. **Figures 3D–F** investigate this range of wider distributions, binning clusters according to their width in each cluster's principal axes, obtained by decomposing each their covariance matrices Σ_{ij} into scaling and rotational components. Principal components are shown in **Figures 3D–F** from most important to least important (namely, Widths 1, 2, and 3), with the height of every bar indicating the number of clusters with a particular width in each principal axis.

Figures 3D–F also show red, green, and blue bars. These bars indicate the number of steady-state clusters with particular widths, namely 1, $\sqrt{2}$, and $\sqrt{3}$. The high cluster count at these particular (Euclidean) distance values is a consequence of the discrete grid-like nature of the agent-based model, which produces steady states such as the ones indicated by the insets in **Figure 3D**, which have widths of 1, $\sqrt{2}$, and $\sqrt{3}$.

Finally, **Figure 3G** shows the frequency of different steady-state clusters, when grouped only regarding their shape, and therefore also aggregating over orientation and centroid position. It confirms what is indicated by the other six panels: the largest fraction of steady-state distributions is point-like, representing all 1,000 agents converging toward the same point in attitude space, a phenomenon which happens for 96% of all model runs, including C matrices with all kinds of symmetry and levels of interdependency between issues. Steady-state distributions two or more states together only take 4% of all runs of the model.

It is important to note that **Figure 3G** is a two-dimensional representation of a three-dimensional model. This is only possible because the frequency of three-dimensional steady states distribution is under 1%, which is comparable to the frequency of other two-dimensional steady states shown in the figure. This is in agreement with **Figure 3F**, which shows that $<1\%$ of all steady states have a non-zero width in their third main axis. In other words: zero-dimensional (point-like) steady states are by far the most common, corresponding to 96% of all model runs, followed by one-dimensional, two-dimensional and three-dimensional steady states, in order of decreasing frequency.

Finally, the reviewer might notice that steady-state distributions such as the bottom left in **Figure 3G** should not be absorbing states under the model with $\beta = 1$. Rather, given enough time, this distribution should converge to the point-like distribution on the top left of **Figure 3G**, which is

an absorbing state. This 1% of all steady-state distributions likely corresponds to runs which are still in their transient state by $t = 15,000$. Preliminary runs of $(M = 10, N = 2)$ and $(M = 10, N = 3)$ show a similar pattern: these system sizes tend to show polarized one-dimensional distributions for timescales longer than 15,000 time steps, only converging to absorbing states after over 50,000 time steps. In their paper, Banisch and Olbrich argue these transient distributions should become more empirically relevant as population sizes grow—we explore this point in more detail in the section 4.

4. DISCUSSION

The aim of this article was to provide a good illustration of the complexity involved in studying an agent-based model of human behavior that is actually guided by social and cognitive psychology. The theoretical details and model choices made by Banisch and Olbrich [35] to model Mäs and Flache's argument communication theory of bi-polarization resulted in a model which is simple to define and to run, but which requires careful analysis, as its outputs are inherently multidimensional and dependent on a number of factors. It is this sort of system which often limits linear and analytical approaches, since the relevant part of the behavior happens at an emergent level. Through a complete enumeration of the $M = 4, N = 2$ and $M = 5, N = 3$ cognitive-evaluative matrices, we find that most runs of the model, for all cognitive-evaluative matrices, move toward a few steady-state distributions. We find that the clusters of steady-state distributions in attitude space corresponding to most runs are often pointwise steady-state distributions, where all agents converge toward the same vector y in attitude space. Steady states composed of two or more attitude states take over approximately 4% of all runs of the model, with distribution with 2 states being the most frequent.

Our analysis of small of Banisch and Olbrich's model for small M and N suggests that the most likely result after many iterations of the model is consensus, and that any deviation from consensus would hardly be described as "polarization." These are, however, small systems: matrices with larger M and N allow for a larger spread of agents in attitude space, which allows for the emergence of distributions polarized along one axis. We observe that in preliminary runs of matrices with $(M = 10, N = 2)$ and $(M = 10, N = 3)$, which display one-dimensional distributions of agents in attitude space for longer than 15,000 time steps, only converging after over 50,000 time steps. This suggests that larger systems should take longer to converge, allowing for the sustained existence of social dynamics within transient population states. Distributions displayed during the transient period should be particularly relevant for larger population sizes, a point also made by Banisch and Araújo when talking more broadly about opinion dynamics models [47].

The main result of this work, beyond producing insights about small systems, is a methodological contribution. As described in more detail in section 2 this multi-level agent-based model does not have any clear output variables, nor a clear aggregation scale, order parameter or measurable outcome.

Its emergent behavior is the product of countless interactions where agents update their beliefs and attitudes, but there is no clear metric assessing when such emergent behavior has happened, or even to tell apart the model transient from its steady state.

This paper introduced a number of approaches to address this problem: in **Figure 2**, after establishing that individual runs of the model for the same cognitive-evaluative matrix should not be averaged without losing significant information, we observe the distribution of agents in attitude space over time, plotting the distance between the agents' centroid over time and the final position of their centroid, as well as looking at the effective number of states of every run. This effective number of states, just like its equivalent measures from other multidimensional models of social behavior, takes an approach from information theory to quantify the diversity of states in the model. With these tools combined, we are able to establish model convergence around $t = 15,000$ steps.

The analysis presented in **Figure 3** presents further methods which can be applied to complex agent-based models: by using a combination of cluster analysis and PCA-like methods to establish the main directions of variation of all the steady states produced by 20 runs of the model for every $M = 5$, $N = 3$. These states were then aggregated in multiple ways, leading to a thorough description of the full spectrum of outputs produced by this model.

The methods presented here open many doors for future research. Firstly, they allow for a more careful exploration of Banisch and Olbrich's model, at system sizes of empirical relevance. Moreover, the full enumeration approach used here might also be ideal—further research is needed to identify the correct ensembles of matrices to represent the mapping between opinions and attitudes. One might also want to consider the

interplay of social network structures and cognitive-evaluative maps, as the separation between beliefs and attitudes might lead to stronger separation between agents in different parts of a network.

Most importantly, this work introduces a scalable way to explore the parameter space of complex agent-based models such as the one studied in this paper. Methods such as the effective number of states or the clustering by steady-state are most appropriate for models which resemble real-life social behavior, particularly the dynamics of beliefs, opinions and attitudes, where emergent phenomena are not static, easily measurable or even clearly defined—and where there usually is no order parameter that identifies different regimes of the model. Here we have introduced not an order parameter, but a set of analysis tools, which can bring more power and clarity to future complex models of social behavior.

DATA AVAILABILITY STATEMENT

The datasets generated for this study are available on request to the corresponding author.

AUTHOR CONTRIBUTIONS

CC designed and performed the research as well as wrote the paper.

ACKNOWLEDGMENTS

The author would like to thank Sven Banisch for helpful discussions on the computational model, and the two reviewers for their very important feedback on previous iterations of this work.

REFERENCES

- Banarjee A. A simple model of herd behaviour. *Q J Econ.* (1992) **107**:797–817.
- Golub B, Jackson MO. Using selection bias to explain the observed structure of internet diffusions. *Proc Natl Acad Sci USA.* (2010) **107**:10833–6. doi: 10.1073/pnas.1000814107
- Golub B, Jackson MO. Naive learning in social networks and the wisdom of crowds. *Am Econ J.* (2010) **2**:112–49. doi: 10.1257/mic.2.1.112
- Axelrod R. The dissemination of culture: a model with local convergence and global polarization. *J Conflict Resol.* (1997) **41**:203–26. doi: 10.1177/0022002797041002001
- Friedkin NE, Johnsen EC. Social influence and opinions. *J Math Sociol.* (1990) **15**:193–206. doi: 10.1080/0022250X.1990.9990069
- Boyd R, Richerson PJ. Why does culture increase human adaptability? *Ethol Sociobiol.* (1995) **16**:125–43. doi: 10.1016/0162-3095(94)00073-G
- Van Benthem J, Van Eijck J, Kooi B. Logics of communication and change. *Inform Comput.* (2006) **204**:1620–62. doi: 10.1016/j.ic.2006.04.006
- Nowak A, Szamrej J, Latané B. From private attitude to public opinion: a dynamic theory of social impact. *Psychol Rev.* (1990) **97**:362. doi: 10.1037/0033-295X.97.3.362
- Clifford P, Sudbury A. A model for spatial conflict. *Biometrika.* (1973) **60**:581–8. doi: 10.1093/biomet/60.3.581
- Holley RA, Liggett TM. Ergodic theorems for weakly interacting infinite systems and the voter model. *Ann Probabil.* (1975) **3**:643–63. doi: 10.1214/aop/1176996306
- Fu F, Wang L. Coevolutionary dynamics of opinions and networks: from diversity to uniformity. *Phys Rev E.* (2008) **78**:016104. doi: 10.1103/PhysRevE.78.016104
- Nardini C, Kozma B, Barrat A. Who's talking first? Consensus or lack thereof in coevolving opinion formation models. *Phys Rev Lett.* (2008) **100**:158701. doi: 10.1103/PhysRevLett.100.158701
- Galam S. Minority opinion spreading in random geometry. *Eur Phys J B-Condens Matter Complex Syst.* (2002) **25**:403–6. doi: 10.1140/epjb/e20020045
- Sznajd-Weron K, Sznajd J. Opinion evolution in closed community. *Int J Modern Phys C.* (2000) **11**:1157–65. doi: 10.1142/S0129183100000936
- Deffuant G, Neau D, Amblard F, Weisbuch G. Mixing beliefs among interacting agents. *Adv Complex Syst.* (2000) **3**:87–98. doi: 10.1142/S0219525900000078
- Burnstein E, Vinokur A. What a person thinks upon learning he has chosen differently from others: nice evidence for the persuasive-arguments explanation of choice shifts. *J Exp Soc Psychol.* (1975) **11**:412–26. doi: 10.1016/0022-1031(75)90045-1
- Burnstein E, Vinokur A. Persuasive argumentation and social comparison as determinants of attitude polarization. *J Exp Soc Psychol.* (1977) **13**:315–32. doi: 10.1016/0022-1031(77)90002-6
- Isenberg DJ. Group polarization: a critical review and meta-analysis. *J Pers Soc Psychol.* (1986) **50**:1141. doi: 10.1037/0022-3514.50.6.1141
- Ewing TN. A study of certain factors involved in changes of opinion. *J Soc Psychol.* (1942) **16**:63–88. doi: 10.1080/00224545.1942.9714105

20. Darke PR, Chaiken S, Bohner G, Einwiller S, Erb HP, Hazlewood JD. Accuracy motivation, consensus information, and the law of large numbers: effects on attitude judgment in the absence of argumentation. *Pers Soc Psychol Bull.* (1998) **24**:1205–15. doi: 10.1177/01461672982411007
21. Kumkale GT, Albarracin D, Seignourel PJ. The effects of source credibility in the presence or absence of prior attitudes: implications for the design of persuasive communication campaigns. *J Appl Soc Psychol.* (2010) **40**:1325–56. doi: 10.1111/j.1559-1816.2010.00620.x
22. Yan Y, Liu J. Effects of media exemplars on the perception of social issues with pre-existing beliefs. *J Mass Commun Q.* (2016) **93**:1026–49. doi: 10.1177/1077699016629374
23. Fiedler K. Meta-cognitive myopia and the dilemmas of inductive-statistical inference. In: Ross B, editor *Psychology of Learning and Motivation, Vol. 57*. Cambridge: Elsevier (2012). p. 1–55.
24. Tversky A, Kahneman D. Belief in the law of small numbers. *Psychol Bull.* (1971) **76**:105. doi: 10.1037/h0031322
25. Goldberg A, Stein SK. Beyond social contagion: associative diffusion and the emergence of cultural variation. *Am Sociol Rev.* (2018) **83**:897–932. doi: 10.1177/0003122418797576
26. Holme P, Newman ME. Nonequilibrium phase transition in the coevolution of networks and opinions. *Phys Rev E.* (2006) **74**:056108. doi: 10.1103/PhysRevE.74.056108
27. Zanette DH, Gil S. Opinion spreading and agent segregation on evolving networks. *Phys D.* (2006) **224**:156–65. doi: 10.1016/j.physd.2006.09.010
28. Bakshy E, Messing S, Adamic LA. Exposure to ideologically diverse news and opinion on Facebook. *Science.* (2015) **348**:1130–2. doi: 10.1126/science.aaa1160
29. Barberá P, Jost JT, Nagler J, Tucker JA, Bonneau R. Tweeting from left to right: is online political communication more than an echo chamber? *Psychol Sci.* (2015) **26**:1531–42. doi: 10.1177/0956797615594620
30. Mäs M, Flache A. Differentiation without distancing. Explaining BI-polarization of opinions without negative influence. *PLoS ONE.* (2013) **8**:e74516. doi: 10.1371/journal.pone.0074516
31. Byrne D. Interpersonal attraction and attitude similarity. *J Abnorm Soc Psychol.* (1961) **62**:713. doi: 10.1037/h0044721
32. Huston TL, Levinger G. Interpersonal attraction and relationships. *Annu Rev Psychol.* (1978) **29**:115–56. doi: 10.1146/annurev.ps.29.020178.000555
33. McPherson M, Smith-Lovin L, Cook JM. Birds of a feather: homophily in social networks. *Annu Rev Sociol.* (2001) **27**:415–44. doi: 10.1146/annurev.soc.27.1.415
34. Wimmer A, Lewis K. Beyond and below racial homophily: ERG models of a friendship network documented on Facebook. *Am J Sociol.* (2010) **116**:583–642. doi: 10.1086/653658
35. Banisch S, Olbrich E. An argument communication model of polarization and ideological alignment. *arXiv [preprint]. arXiv:180906134* (2018).
36. Fishbein M, Raven BH. The AB scales: an operational definition of belief and attitude. *Hum Relat.* (1962) **15**:35–44. doi: 10.1177/001872676201500104
37. Fishbein M. An investigation of the relationships between beliefs about an object and the attitude toward that object. *Hum Relat.* (1963) **16**:233–9. doi: 10.1177/001872676301600302
38. Ajzen I. Nature and operation of attitudes. *Annu Rev Psychol.* (2001) **52**:27–58. doi: 10.1146/annurev.psych.52.1.27
39. Rosa H. *Resonance: A Sociology of Our Relationship to the World*. Hoboken, NJ: John Wiley & Sons (2019).
40. Dalege J, Borsboom D, van Harreveld F, van den Berg H, Conner M, van der Maas HL. Toward a formalized account of attitudes: the Causal Attitude Network (CAN) model. *Psychol Rev.* (2016) **123**:2–22. doi: 10.1037/a0039802
41. Krause U, Stöckler M. *Modellierung und Simulation von Dynamiken mit Vielen Interagierenden Akteuren*. Bremen: Modus Universität (1997).
42. Hegselmann R, Krause U. Opinion dynamics and bounded confidence models, analysis, and simulation. *J Artif Soc Soc Simul.* (2002) **5**. Available online at: <http://jasss.soc.surrey.ac.uk/5/3/2.html>
43. Deffuant G, Amblard F, Weisbuch G, Faure T. How can extremism prevail? A study based on the relative agreement interaction model. *J Artif Soc Soc Simul.* (2002) **5**:1. Available online at: jasss.soc.surrey.ac.uk/5/4/1.html
44. Weisbuch G, Deffuant G, Amblard F, Nadal JP. Meet, discuss, and segregate! *Complexity.* (2002) **7**:55–63. doi: 10.1002/cplx.10031
45. Stauffer D, De Oliveira SMM, De Oliveira PMC, de Sá Martins JS. *Biology, Sociology, Geology by Computational Physicists*. Amsterdam: Elsevier (2006).
46. Lorenz J. Continuous opinion dynamics under bounded confidence: a survey. *Int J Modern Phys C.* (2007) **18**:1819–38. doi: 10.1142/S0129183107011789
47. Banisch S, Araújo T, Louçã J. Opinion dynamics and communication networks. *Adv Complex Syst.* (2010) **13**:95–111. doi: 10.1142/S0219525910002438
48. Acemoglu D, Ozdaglar A. Opinion dynamics and learning in social networks. *Dyn Games Appl.* (2011) **1**:3–49. doi: 10.1007/s13235-010-0004-1
49. Shannon CE. A mathematical theory of communication. *Bell Syst Tech J.* (1948) **27**:379–423. doi: 10.1002/j.1538-7305.1948.tb01338.x
50. Jost L. Entropy and diversity. *Oikos.* (2006) **113**:363–75. doi: 10.1111/j.2006.0030-1299.14714.x
51. Laakso M, Taagepera R. “Effective” number of parties: a measure with application to West Europe. *Compar Polit Stud.* (1979) **12**:3–27. doi: 10.1177/001041407901200101
52. Camargo CQ, Hale SA, John P, Margetts HZ. Volatility in the issue attention economy. *arXiv [preprint]. arXiv:180809037.* (2018).

Conflict of Interest: The author declares that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2020 Camargo. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.



How Could Future AI Help Tackle Global Complex Problems?

Anne-Marie Grisogono*

College of Science and Engineering, Flinders University, Adelaide, SA, Australia

How does AI need to evolve in order to better support more effective decision-making in managing the many complex problems we face at every scale, from global climate change, collapsing ecosystems, international conflicts and extremism, through to all the dimensions of public policy, economics, and governance that affect human well-being? Research in complex decision-making at an individual human level (understanding of what constitutes more, and less, effective decision-making behaviors, and in particular the many pathways to failures in dealing with complex problems), informs a discussion about the potential for AI to aid in mitigating those failures and enabling a more robust and adaptive (and therefore more effective) decision-making framework, calling for AI to move well-beyond the current envelope of competencies.

OPEN ACCESS

Keywords: AI decision support, complex decisions, human limitations, wicked problems, interface design

Edited by:

Carlos Gershenson,
National Autonomous University of
Mexico, Mexico

Reviewed by:

Melanie E. Moses,
University of New Mexico,
United States
Caleb Rascon,
National Autonomous University of
Mexico, Mexico

*Correspondence:

Anne-Marie Grisogono
annemarie.grisogono@flinders.edu.au

Specialty section:

This article was submitted to
Computational Intelligence in
Robotics,
a section of the journal
Frontiers in Robotics and AI

Received: 03 November 2019

Accepted: 20 March 2020

Published: 21 April 2020

Citation:

Grisogono A-M (2020) How Could
Future AI Help Tackle Global Complex
Problems? *Front. Robot. AI* 7:50.
doi: 10.3389/frobt.2020.00050

INTRODUCTION

Human intelligence rests on billions of years of evolution from the earliest origins of life, and despite its undeniably unique nature within the biosphere, and the apparent gulf that distinguishes the human species from all others, it should nevertheless be seen as an extremum within a continuum. The unifying feature of all natural intelligence systems is that they have evolved under strong selection pressures to solve the problems of surviving and thriving sufficiently to reproduce better than their competitors. Unlike the evolution of faster speed, sharper teeth, more efficient energy harvesting and utilization, or better camouflage, all of which improve physical capabilities, the evolution of intelligence enables better choices to be made as to how and when to employ those capabilities, by processing relevant sensed and stored information. If the environment is challenging enough, whether through the prevalence of threats, the scarcity of necessary resources or through intense competition for them, then there is a high fitness pay-off for evolving both the necessary physical characteristics for sensing, processing and storing the relevant information, and the intelligence to exploit them.

From this perspective we can define intelligence as the ability to produce effective responses or courses of action that are solutions to complex problems—in other words, problems that are unlikely to be solved by random trial and error, and that therefore require the abilities to make finer and finer distinctions between more and more combinations of relevant factors and to process them so as to generate a good enough solution. Obviously this becomes more difficult as the number of possible choices increases, and as the number of relevant factors and the consequence pathways multiply. Thus complexity in the ecosystem environment generates selection pressure for effective adaptive responses to the complexity.

One possible adaptive strategy is to find niches to specialize for, within which the complexity is reduced. The opposite strategy is to improve the ability to cope with the complexity by evolving increased intelligence at an individual level, or collective intelligence through various types of

cooperative or mutualistic relationships. Either way, increased intelligence in one species will generally increase the complexity of the problems they pose for both other species in the shared ecosystem environment, and for their own conspecifics, driving yet further rounds of adaptations. Even when cooperative interactions evolve to deal with problems that are more complex than an individual can cope with, the shared benefits come with a further complexity cost in maintaining the cooperative relationships and policing for cheats (Nowak, 2006).

This ratcheting dynamic of increasing intelligence and increasing complexity continues as long as two conditions are met: further increases in sensing and processing are sufficiently accessible to the evolutionary process, and the selection pressure is sufficient to drive it. Either condition can fail. Thus generally a plateau of dynamic equilibrium is reached. But it is also possible that under the right conditions, which we will return to below, the ratcheting of both complexity and intelligence may continue and accelerate.

Artificial intelligence on the other hand, has not evolved through natural selection, but rather owes its genesis to human intelligence (at least on this planet), which has a number of important implications that have colored its trajectory so far. But to contemplate its possible futures and ours, this paper argues the need to re-examine the relationship between human and machine within a much broader context. In particular, we need to understand both the strengths and the limitations of human intelligence, consider what our most pressing issues are and what kinds of advances in AI would be most useful in helping us to navigate those complex problems in the near to mid-term. At the same time we need to be mindful of the risks, not only in the nearer term but also those that may only materialize as longer term consequences, and address how these may be averted or mitigated.

AI AND THE LIMITATIONS OF HUMAN INTELLIGENCE

It was natural for the pioneers of AI to choose human cognitive abilities such as playing chess or Go, navigating obstacles, or recognizing and interpreting written and spoken language, as the yardsticks by which to measure early progress in AI capabilities, not only because they were so far beyond what could be simulated at the time, but also perhaps, because we felt so impressed with our own dazzling cognitive strengths. But now that many of these and other quintessentially human examples of intelligence are being relegated to the growing list of tasks at which AI can surpass human performance, we need to step back and acknowledge that human intelligence is not the pinnacle of what can be achieved.

Just as the Copernican revolution and later astronomical discoveries dislodged us from the center of the universe and pushed us into orbiting a minor star in an undistinguished galaxy, and Darwinism pushed us from the pre-eminent position we had assumed over all life forms into just a twig of the evolutionary tree of life, the current and recent sweep of advances in understanding of neuroscience, cognition, behavioral science, evolutionary psychology and related fields call for yet another

round of humbling re-appraisal of where we fit in the grand scheme of things.

Taking the concept of intelligence as the ability to produce effective solutions to complex problems by processing relevant sensed and stored information, it is evident that human intelligence and ingenuity have led to immense progress in producing solutions for many of the pressing problems of past generations, such as higher living standards, longer life expectancy, better education and working conditions. But it is equally evident that the transformations they have wrought in human society and in the planetary environment include many harmful unintended consequences, and that the benefits themselves are not equitably distributed and have often masked unexpected downsides.

We are now confronting a complex network of interdependent global problems which we seem increasingly incapable of dealing with effectively at either the national or international levels, and arguably it is the very successes of human intelligence that have ratcheted the complexity of the challenges we face to a level that unaided human intelligence is now unable to cope with.

This was recognized as long ago as 1973 in a remarkably prescient paper (Rittel and Webber, 1973) in which the authors coined the term “wicked problems” (as opposed to benign problems which are tractable) and laid out ten hallmarks¹ characterizing them, together with a very clear analysis of their roots in complexity. Their inability to lay out an equally clear prescription for the resolution of such wicked problems signaled that a tipping point had indeed been reached where our limitations had now outstripped our cleverness.

What has changed in the intervening decades? While the scale and urgency of the global problems we face have certainly intensified, what we have since learned in the germane fields of complexity science, evolutionary psychology, brain and behavioral science, and artificial intelligence, suggests that we may be close to another tipping point where we could possibly drive the emergence of advanced artificial intelligence systems that can effectively support human decision-making in managing such problems, by a combination of mitigating human fallibilities and complementing human shortcomings.

At this point the reader may be wondering why there should be a human in the decision process at all if we have indeed overstepped our domain of competence. There are possibly three reasons.

Firstly, even if there does come a day when AI systems are judged able to take over the management of complex issues without human control, such a judgment would imply that humans have confidence in those systems, and such confidence can only be developed through a transition period of human and machine working together, learning the strengths and

¹Briefly the ten hallmarks are: no definitive formulation; no stopping rule; solutions are not true-or-false, but better or worse; no immediate or ultimate test of a solution; every solution attempt is a “one-shot operation”; no well-described set of potential solutions or permissible operations; essentially unique; can be considered a symptom of another problem; many possible explanations; and the decision-maker has “no right to be wrong” because of the gravity of the consequences.

limits of each other's capabilities, and evolving better ways to arrive at good decisions, through evaluating and learning from the consequences of those decisions. Secondly, there is the perennial issue of expert knowledge elicitation. Despite all their human failings, there is surely a vast, unquantifiable reservoir of relevant experiential implicit knowledge, and hopefully wisdom, in the cohorts of public officials, managers and analysts who currently strive to deal with these spiraling problems. If an AI system is to eventually run things without them, it had better somehow absorb what they know that cannot be itemized in databases - which links to the third reason: will people really want to be excluded from managing their societies and enterprises? The answer that might emerge in the future when the question actually becomes pertinent is impossible to predict today. But we have enough reasons to proceed on the assumption that the next steps will involve advanced AI support for human decision-makers.

To propose a set of desiderata for the advances in AI that are needed we now turn to what we have learned about the specific limitations that plague human decision-makers in complex problems. We can break this down into two parts: the aspects of complex problems that we find so difficult, and what it is about our brains that limits our ability to cope with those aspects.

Sources of Difficulty in Complexity

Interdependence is a defining feature of complexity and has many challenging and interesting consequences. In particular, the network of interdependencies between different elements of the problem means that it cannot be successfully treated by dividing it into sub-problems that can be handled separately. Any attempt to do that creates more problems than it solves because of the interactions between the partial solutions.

Dynamical processes driving development of the situation often involve many positive and negative feedbacks, thus amplifying and suppressing different aspects of the situation, and resulting in highly **non-linear dynamics**. This means that relying on linear extrapolation of current conditions can lead to serious errors.

There is no natural boundary that completely isolates a complex problem from the context it is embedded in. There is always some traffic of information, resources, and agents in and out of the situation which can bring about unexpected changes, and therefore the **context cannot be excluded** from attention.

Complex problems exist at **multiple scales**, with different agents, behaviors and properties at each, but with interactions between scales. This includes both emergence, the appearance of complex structure and dynamics at larger scales as a result of smaller-scale phenomena, and its converse, top-down causation, whereby events or properties at a larger scale can alter what is happening at the smaller scales. In general, all the scales are important, there is no single "right" scale at which to act.

Interdependence implies multiple interacting causal and influence pathways leading to, and fanning out from, any event or property, so simple causality (one cause—one effect), or linear causal chains will not hold in general. Yet much of our cultural conditioning is predicated on a naïve view of linear causal chains, such as finding "the cause" of an effect, or "the person" to be

held responsible for something, or "the cure" for a problem. Focusing on singular or primary causes makes it more difficult to intervene effectively in complex systems and produce desired outcomes without attendant undesired ones—so-called "side-effects" or **unintended consequences**. Effective decision making requires the ability to develop sufficient understanding of the causal and influence network to engage with it effectively, neither oversimplifying it, nor becoming overwhelmed with unnecessary levels of detail.

Furthermore, such networks of interactions between contributing factors can produce **emergent behaviors** which are not readily attributable or intuitively anticipatable or comprehensible, implying unknown risks and unrecognized opportunities.

There are generally **multiple interdependent goals** in a complex problem, both positive and negative, poorly framed, often unrealistic or conflicted, vague or not explicitly stated, and stakeholders will often disagree on the weights to place on the different goals, or change their minds. Achieving sufficient high level goal clarity to develop concrete goals for action is in itself a complex problem.

Complex situations generally contain many **adaptive agents** with complex relationships and shifting allegiances, and new behaviors and features continually arise. This means that approaches that worked in the past may no longer work, interventions that frustrate the intents of some agents will often simply stimulate them to find new ways to achieve them, and opportunities created by the inevitable new vulnerabilities that interventions create will be rapidly identified and exploited.

Many important aspects of complex problems are hidden, so there is **inevitable uncertainty** as to how the events and properties that are observable, are linked through causal and influence pathways, and therefore many hypotheses about them are possible. These cannot be easily distinguished based on the available evidence.

Limitations of the Human Brain

The brief overview above reveals some of the cognitive abilities that are essential for successful tackling of complex problems. One immediate conclusion that can be drawn is that there is a massive requirement for cognitive bandwidth—not only to keep all the relevant aspects at all the relevant scales in mind as one seeks to understand the nature of the problem and what may be possible to do, but even more challenging, to incorporate appropriate non-linear dynamics as trajectories in time are explored. Given the well-known limitations of human working memory, short-term memory and attention span, this is an obvious area for advanced AI support to target.

But there is a more fundamental problem that needs to be addressed first: how to acquire the necessary relevant information about the composition, structure and dynamics of the complex problem and its context at all the necessary scales, and revise and update it as it evolves. This requires a stance of continuous learning, i.e., simultaneous sensing, testing, learning and updating across all the dimensions and scales of the problem, and the ability to discover and access relevant sources of information. At their best, humans are okay at this, up to a

point, but not at the sheer scale and tempo of what is required in real world complex problems which refuse to stand still while we catch up.

Moreover, there are both physiological factors such as the impacts of stress, fatigue and anxiety on cognitive performance, and particular features of the human brain, legacies of our evolutionary history, which compound the difficulties.

Because the human brain evolved to deal with the problems of surviving and thriving that our ancestors faced, modern humans are still equipped with the same heuristics, behavioral tendencies and biases that worked well enough in the distant past. These hardwired shortcuts based on rules of thumb, operating automatically below conscious awareness and so permitting very rapid adaptive responses to various simple conditions, enabled them to cope with the level of complexity that existed then—keeping track of a hundred or so individuals and their interactions, intents, and histories (Dunbar, 1992). But features relying on approximations that held true for dealing with common problems in past environments can morph into risky bugs in today's highly interconnected and rapidly evolving complex situations (Kahneman, 2002).

To understand how all these factors interact to limit human competence in managing complex problems, and what opportunities might exist for mitigating them through advanced AI systems, we now review some key findings from relevant research.

In particular we are interested in learning about the nature of human decision-making in the context of attempting to manage an ongoing situation which is sufficiently protracted² and complex to defeat most, but not all³, decision-makers. Drawing useful conclusions about the detailed decision-making behaviors that tend to either sow the seeds of later catastrophes, or build a basis for sustained success, calls for an extensive body of empirical data from many diverse human subjects making complex decisions in controllable and repeatable complex situations. Clearly this is a tall ask, so not surprisingly, the field is sparse. However, one such research program (Dörner, 1995; Evans et al., 2011; Dörner and Gerdes, 2012; Dörner and Güss, 2013; Donovan et al., 2015), which has produced important insights about how successful and unsuccessful decision-making behaviors differ, stands out in having also addressed the underlying neurocognitive and affective processes that conspire to make it very difficult for human decision-makers to maintain the more successful behaviors, and to avoid falling into a vicious cycle of less effective behaviors.

In brief, through years of experimentation with human subjects attempting to achieve complex goals in computer-based micro-worlds with complex underlying dynamics, the specific

²Managing complex situations involves many decisions over an extended period, with the consequences of earlier ones impacting on the necessity or possibility of later ones, and affecting the trajectory of the situation. To come to grips with how decision-making behaviors shape outcomes it is important to conduct experiments for a long enough period to allow these consequences to develop and confront the decision-maker.

³In order to learn what decision-making behaviors are more effective, the degree of complexity of the experimental environment has to be tuned to the edge of human competence so that data can also be gathered about what does work.

decision-making behaviors⁴ that differentiated a small minority of subjects who achieved acceptable outcomes in the longer term, from the majority who failed to do so, were identified. Results indicated that most subjects could score some quick wins early in the game, but as the unintended consequences of their actions developed and confronted them, and their attempts to deal with them created further problems, the performance of the overwhelming majority (~90%) quickly deteriorated, pushing their micro-worlds into catastrophic or chronic failure.

As would be expected, their detailed behaviors reproduced many well-documented findings about the cognitive traps posed by human heuristics and biases. Low ambiguity tolerance was found to be a significant factor in precipitating the behavior of prematurely jumping to conclusions about the problem and what was to be done about it, when faced with situational uncertainty, ambiguity and pressure to achieve high-level goals. The chosen (usually ineffective) course of action was then defended and persevered with through a combination of confirmation bias (Nickerson, 1998), commitment bias (Staw, 1997), and loss aversion (Kahneman and Tversky, 1979), in spite of available contradictory evidence. The unfolding disaster was compounded by a number of other reasoning shortcomings such as difficulties in steering processes with long latencies and in projecting cumulative and non-linear processes (Serman, 1989). Overall they had poor situation understanding, were likely to focus on symptoms rather than causal factors, were prone to a number of dysfunctional behavior patterns, and attributed their failures to external causes rather than learning from them and taking responsibility for the outcomes they produced.

By contrast, the remaining ten percent who eventually found ways to stabilize their micro-world, showed systematic differences in their decision-making behaviors and were able to counter the same innate tendencies by taking what amounts to an adaptive approach, developing a conceptual model of the situation, and a stratagem based on causal factors, seeking to learn from unexpected outcomes, and constantly challenging their own thinking and views. Most importantly, they displayed a higher degree of ambiguity tolerance than the unsuccessful majority.

These findings are particularly significant here because most of the *individual* human decision-making literature has concentrated on how complex decision-making fails, not on how it succeeds. However, insights from research into successful *organizational* decision-making in complex environments (Collins, 2001; Weick and Sutcliffe, 2001), do corroborate the importance of taking an adaptive approach.

In summary, analysis of the effective decision behaviors offers important insights into what is needed, in both human capabilities and AI support, to deal with even higher levels of complexity beyond current human competence. There are two complementary aspects here—put simply: how to avoid pitfalls (what not to do), and how to adopt more successful approaches (what to do instead).

⁴The behaviors were grouped in five categories: goal decomposition; collecting and organizing information; projection and planning; decision and execution; and meta-cognition.

It is not difficult to understand how the decision making behaviors associated with the majority contributed to their lack of success, nor how those of the rest enabled them to develop sufficient conceptual and practical understanding to manage and guide the situation to an acceptable regime. Indeed if the two lists of behaviors are presented to an audience, everyone can readily identify which list leads to successful outcomes and which leads to failure. Yet if those same individuals are placed in the micro-world hot seat, 90% of them will display the very behaviors they just identified as likely to be unsuccessful. This implies that the displayed behaviors are not the result of conscious rational choice, but are driven to some extent by unconscious processes.

This observation informed development of a theoretical model (Dörner and Gerdes, 2012; Dörner and Güss, 2013) incorporating both cognitive and neurophysiological processes to explain the observed data. In brief, the model postulates two basic psychological drives that are particularly relevant to complex decision making, a need for certainty and a need for competence. These are pictured metaphorically as tanks which can be topped up by signals of certainty (one's expectations being met) and signals of competence (one's actions producing desired outcomes), and drained by their opposites—surprises and unsuccessful actions. The difference between the current level and the set point of a tank creates a powerful unconscious need, stimulating some behavioral tendencies and suppressing others, and impacting on cognitive functions through stimulation of physiological stress. If both levels are sufficient the result is motivation to explore, reflect, seek information and take risky action if necessary—all necessary components of effective decision making behavior. But if the levels get too low the individual becomes anxious and is instead driven to flee, look for reassurance from others, seek only information that confirms his existing views so as to top up his dangerously low senses of certainty and competence, and deny or marginalize any tank-draining contradictory information. The impacts of stress on cognitive functions reinforce these tendencies when the levels are too low by reducing abilities to concentrate, sustain a course of action, and recall relevant knowledge.

Individuals whose tanks are low therefore find it difficult to sustain the decision-making behaviors associated with success, and are likely to act in ways that generate further draining signals, digging themselves deeper into a vicious cycle of failure. We can now understand the 90:10 ratio, as the competing attractors are not symmetric—the vicious cycle of the less effective decision behaviors is self-reinforcing and robust, while the virtuous cycle of success is more fragile because one's actions are not the sole determinant of outcomes in a complex situation, so even the best decision-makers will sometimes find their tanks getting depleted, and therefore have difficulty sustaining the more effective decision making behaviors.

Further research has demonstrated that the more effective decision making behaviors are trainable to some extent, but because they entail changing meta-cognitive habits they require considerable practice, reinforcement and ongoing support (Evans et al., 2011; Grisogono and Radenovic, 2011; Donovan et al., 2015). However, the scope for significant enhancement of unaided human complex decision making competence is

limited—not only in the level of competence achievable, but also and more importantly, in the degree of complexity that can be managed.

Meanwhile, the requirements for increased competence, and the inexorable rise in degree of complexity to be managed, continue to grow.

How Could AI Help?

Recent AI advances such as deep learning and generative adversarial networks have demonstrated impressive results in many domains—superhuman precision in classification tasks, beating human world champions in Go, and generation of images that are hard for humans to discriminate from reality, to name a few.

But what are the prospects for advances in AI to deliver the kind of decision support capability that is needed by those charged with managing the most challenging, indeed wicked, problems? And can those advances be achieved by research that continues to set goals based on beating human performance, or on fooling human discrimination?

Despite its successes, the best examples of AI are still very specialized applications that focus on well-defined domains, and that generally require a vast amount of training data to achieve their high performance. Such applications can certainly be components of an AI decision support system for managing very complex problems, but the factors discussed in the two previous sections imply that much more is needed: not just depth in narrow aspects, but breadth of scope by connecting the necessary components so as to create a virtual environment which is a sufficiently valid model of the problem and its context, and in which decision-makers can safely explore and test options for robustness and effectiveness, while being supported in maintaining effective decision making behaviors and resisting the less effective ones. The following section develops a more detailed set of desiderata for such an AI support system.

The resurgence of interest in Artificial General Intelligence seems a promising avenue for the kinds of advances that are needed, but it is telling that AGI is most often explicitly pursued through the lens of the touted general intelligence that humans possess (Adams et al., 2012), in other words still focusing on what we believe we are good at, rather than exploring the most critical parts of the very much larger space of what we are not good enough at. But is human intelligence truly general? The claim rests principally on our ability to learn, and this is certainly a core requirement for future intelligent systems. But we should also acknowledge that the human brain is the product of our particular evolutionary history and sports the evidence of its contingencies in many kluges, biases and peculiarities (Marcus, 2009). It would be reasonable to suppose that other more efficient, more general, more powerful and less flawed designs are possible.

Obviously there is still an immense amount to be learned about how human intelligence actually works and how the detailed structure and architecture of the brain produces it. And there will certainly be many insights that can be implemented in novel AI developments—for example the recent breakthroughs in understanding the workings of the neocortex (Hawkins et al.,

2019), and the comprehensive program to develop cognitive models of how humans build compositionally structured causal models of the world grounded in their capacities for intuitive physics and intuitive psychology, so as to apply them to development of advanced AI systems (Tenenbaum et al., 2011; Lake et al., 2017). However, there is also an argument to be made that relying too much on guiding further development of AI on what is known about human intelligence, risks reproducing some of its limitations, or at least misses opportunities to deliberately and specifically mitigate them so as to extend and complement human capability.

DESIDERATA FOR AN AI DECISION SUPPORT SYSTEM FOR COMPLEX PROBLEMS

The preceding discussion suggests an AI decision support system with three functional areas: an interface through which the human decision-maker interacts with it, the AI core generating and operating on a virtual conceptual model, and an interface to the outside world through which the AI core can grow its capability. Since the future system envisaged here is well beyond what is currently possible, its design can only be sketched out conceptually. The following two subsections offers some high level desiderata for the interfaces and the core, based on a hypothetical use case: in the light of the research insights presented in the preceding section, what would be most useful to a well-intentioned human decision-maker faced with very complex situations to manage?

A third subsection raises some of the ethical issues that must be addressed if such a system is able to be built.

Interfaces to the Human Decision-Maker and the Outside World

The decision-maker⁵ needs to be able to give the system some initial direction about the problem, its scope, context, and goals and then develop them through dialogue, with intuitive visualizations presented by the interface to anchor and stimulate his participation. As these take shape the dialogue should extend to exploration of possible actions and their consequences, the development of courses of action, the building of necessary support from stakeholders and eventually monitoring the implementation of decisions made, and revising all above as more is learned and as the situation evolves.

The way that these are presented should support human understanding of the emerging conceptual model of the problem and its context, implying an appropriate level of coarse-graining in terms of intuitively comprehensible parameters. In particular, the interface should expose both explicit and implicit assumptions in the conceptual model, and possible levers of action and their consequences, both in the intended pathways and in other pathways that may be stimulated, together with

estimates of the degree of uncertainty and the risks resulting from the consequent ranges of possible outcomes.

To reduce risks and further develop the conceptual model, the ongoing dialogue between the human and the interface should be able to launch searches for more data, initiate probing actions, and pose and explore “what if?” and “how could?” questions.

Conflicts and trade-offs also need to be identified – both those that must be explicitly managed, such as the balance between long-term and short-term outcomes, competing interests between different agents, and conflicts between espoused values and/or principles, and those that are actually false dichotomies which should be resolved by supporting exploration of integrative solutions in their place.

Most importantly, to enable the necessary adaptive approach, the interface must not only continuously evolve the conceptual model, but also in parallel prompt and support a process of continuous co-evolution of the goals, data collection plan, and both the structure and implementation of the strategy.

The decision-maker needs to have confidence that the system is in fact presenting accurate and comprehensive information and making judgments in accord with a transparent and agreed set of goals, values and principles. This implies additional requirements with respect to visibility of the goals, values and principles on which it is operating, flagging of uncertainties and assumptions, and where possible testing them, and demonstration that it is using its searching and learning resources to improve its conceptual model so as to reduce risks and uncertainties, in other words, actively subjecting its critical aspects to severe testing, and generating an audit trail for decisions made in relation to every complex issue.

These considerations imply that both the interface and the conceptual model behind it must be open systems that permit evolution of the vocabulary of the interface and the semantic map to the ontology of the model.

Since the interface is also the locus of the metacognitive support that the system can provide to the decision-maker, its design must be informed by an understanding of human limitations and shortcomings.

In particular, and building on, but going beyond the currently established principles of human computer interaction, for which a vast literature exists⁶, the interface design should scaffold a human decision-maker who seeks to overcome the specific difficulties and obstacles discussed in the previous section. For example, the interface could monitor for the influence of unconscious biases in the decision-maker’s actions, such as confirmation bias, framing and recency biases, loss/gain asymmetry and so on, flagging them for conscious attention and offering options for reducing them. It could also reduce anxiety stemming from ambiguity, by demonstrating that an effective risk management strategy is in place (i.e., that indicators of emerging risks are being monitored and averting or mitigating action plans are ready to be triggered), and anxiety stemming from information overload, by effective partitioning of the

⁵For simplicity here we assume a single decision-maker, recognizing that in a real world problem situation there will be many involved and that will necessitate further support requirements.

⁶See for example the scope of CHI ’19- Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems <http://st.sigchi.org/publications/toc/chi-2019.html>

information assimilation workload. Shouldering the workload of maintaining as many alternative working hypotheses as necessary, and exposing and testing the implicit assumptions in each of them, would assist in reducing the danger of premature convergence to a narrow singular view of the problem and hence selecting an inadequate strategy.

Overall what is being described here is a cooperative system where learning and adaptation occurs at the levels of both the human and the AI support system. Importantly, it also occurs at the level of the combined system—the interface supporting the decision-maker's learning by setting the example of its own learning behavior, in particular by continuously making predictions based on its current conceptual model, monitoring for the real world outcomes and revising its models in the light of what has been learned, and the human decision-maker being willing to expose their reasoning and ideas and subject them to analysis in their dialogue with the AI support system.

Of course these hypothetical examples are illustrative rather than prescriptive, and certainly not comprehensive, but they can serve as an adequate starting place for an iterative and continuously learning design process for the interface. Future research will no doubt surface many further opportunities to enhance both human decision-making performance, and the decision performance of the combined human plus AI support system.

Irrespective of how this design research agenda might evolve in detail in the future, one important consequence that seems inevitable is that the system's interface to the outside world must be able to autonomously engage in all relevant aspects as a trusted partner or agent of the decision-maker. Therefore, to enable the scale, tempo and depth of testing and learning that is called for in dealing with multi-faceted and open-ended complex problems, the interface to the outside world must be essentially unfettered and support multiple simultaneous high bandwidth interactions, as well as robust and secure. This point will have repercussions in discussing ethical concerns.

AI Core

These requirements imply that the AI core needs the ability to develop situational models of the complex problems to be managed, and as much of their context as necessary, and to evolve them in a real time loop through predictive processing (Clark, 2015) and updating, i.e., by monitoring relevant developments, using the current version of the model to predict expected consequences, comparing predictions to actual outcomes, and hence updating the models as a result of what is learned. This means that the models must be open systems so that their structures and composition can change as more is learned, and as the situation itself changes over time.

The models need to exist at multiple scales—from coarse resolution to as fine a level as is required to model the relevant entities and events (whether by bottom-up models or by machine learning from data), and include all the dimensions relevant to the necessary scales of representation and all significant outcome variables, all accessible levers of influence that could be exercised, all the causal and influence pathways that may lead to significant consequences, the causal and influence relationships

between entities and events, within and across scales, and their time dependence.

Including all the significant outcome variables implies a detailed representation of how success and failure of the complex problem will be judged, as well as intermediate outcomes and indicators that signal which consequence pathways are activated.

Situation models with such wide scope will necessarily be hybrid models, containing many detailed components, plus representations of the interactions and interdependencies between the components. To support zooming between scales, the core will need the ability to extract human comprehensible coarse-grained models⁷ from the more detailed models, whether data driven ML models or bottom-up micro-parameter based models.

To deal with complex problems at the “wicked” end of the complexity scale the core will need to be able to model humans who are stakeholders or actors in the situation, so that their responses to interventions or external events can be anticipated, and combinations of incentives and compensatory measures can be discovered that have a chance of fostering enough consensus for effective action to be taken. These models will also have to be learned by predictive processing, and continuously tested and updated.

In particular, the core will need to develop very good models of the human decision-makers which it is supporting, so that it can learn to interact with them in a way that they will value and trust.

In summary, besides the requisite models, the AI core needs a number of intelligent functions to enable all the operations implied by the considerations above, and the ability to evolve these as well in the light of its experience and interactions with the decision-maker in order to improve its capability.

Ethical Issues

If such an intelligent support system is ever built it will be extremely useful and powerful. How could misuse be prevented? This is a serious question which must be addressed at the earliest stages of development. Internationally agreed guidelines⁸ and regulation, and a set of safety standards to be met, together with public transparency of the setup and use of any such system would at least make it possible to monitor the known systems. Detecting covert systems is more challenging and may need to be part of an overall cybersecurity capability, along with ensuring security from malicious manipulation.

The requirement for the system to be an autonomous agent with broad unlimited access to the world for learning and testing purposes, will raise particular ethical issues not only with respect to privacy, but also with respect to the commonly expressed fear that as AI becomes more intelligent and powerful, it will become harder, if not impossible, to continue to exert human control over it. The need to allow it to become more autonomous and intelligent and situated in the real world in order to be sufficiently

⁷This has proven difficult so far but recent work in this paper, Mattingly et al. (2018) and references therein may provide a breakthrough, not only in generating a human-comprehensible coarse-grained model, but importantly in identifying the few “stiff” parameter combinations that characterize its emergent macro level properties.

⁸See Jobin et al. (2019) for a recent overview.

effective is at odds with one proposed safeguard measure—that of strictly limiting its access to other real world systems. Therefore, it will be essential to develop other approaches to ensure that it continues to serve human needs and interests. However, defining those needs and interests will be a challenging and controversial wicked problem in itself. All these considerations point to the importance and urgency of addressing the ethical issues at the earliest stages.

Transparency is a powerful aid to addressing some of the ethical issues with AI supported decisions that may have adverse impacts on individuals or groups. For example prejudicial bias introduced into machine learning systems through training data could be exposed through triangulation with independent data sources. Similarly, exposing the assumptions that are made in the conceptual model, together with the efforts that have been made to test them, and whatever evidences are available to support or refute them, would help ethics watchdogs do their job.

DISCUSSION

It was noted in the Introduction section that the mutual ratcheting of complexity and intelligence did not necessarily terminate in a plateau of dynamic equilibrium. Under the right conditions it could continue and accelerate.

The right conditions are that selection pressure for intelligence remains strong and that the evolutionary process is able to generate further improvements in intelligence.

This describes where we are today. We desperately need more powerful intelligence to navigate the perilous waters we find ourselves in, and we have spawned completely new channels of creating and evolving intelligence beyond those afforded purely by our own biology. And both processes are arguably accelerating. Therefore, we do not have the option of turning back.

But it does raise another serious ethical question: where will the ratcheting dynamic of complexity and intelligence lead us? Will AI-aided resolution (or at least diminution) of tomorrow's most serious global problems generate even more disastrously wicked problems in a chain of escalation that rapidly drives humans to irrelevance?

While we cannot rule out worst case fears, the preceding discussions suggest two considerations that give grounds for cautious optimism.

Firstly the “wickedness” of wicked problems is in large part due to the shrinking of the viable option space as more agents with diverse priorities acquire a veto stake in the decision process and so need to be simultaneously satisfied. But a future AI support system could ameliorate this problem, through its capabilities to model the different agents and to devise strategies to win them over—as has already been demonstrated several times recently in the manipulation of voter opinions and preferences (Burkell and Regan, 2019). Of course this also raises ethical concerns and there would need to be a code of conduct agreed that provided transparency and guidelines as to what was acceptable.

Secondly, the ratcheting of complexity observed so far has largely been driven by short-sighted “fixes” of perceived problems, without much consideration of longer term and wider scope consequences—hence inadvertently creating further problems. This is intrinsically the case in natural evolutionary processes, and also very much the case with human decision-makers due to their limited cognitive bandwidth. (A good example is the rapid evolution of mines that are more lethal and harder to detect being driven by researching and fielding better vehicle protection and mine detection systems.) Again, a future AI support system could potentially reduce the pace of ratcheting, by anticipating longer term and wider scope consequences, factoring them in to the evaluation of strategy options, and where necessary actively reducing unwanted consequences with further supplementary actions.

Of course there is also the possibility that the tide does not turn, but rather continues to pose growing threat levels. But then what choice do we have? The immediate global and national problems facing us are urgent and we need all the help we can get. If we decline the opportunity to develop such systems, we will in any case face escalating problems, which might now include opponents and vested interests armed with the very capabilities we declined.

This suggests that it is time to shift the balance of investment in AI research and development away from competing with humans and toward creating new cooperative partnerships with them, to extend and buttress our joint capability to manage the rafts of wicked problems that threaten us. It will involve developing many new aspects of AI capability, but every new capability we create will help generate the next. We are rushing into a future that we can barely imagine, but we need to look ahead with as much clarity as we can muster, embrace the present opportunity we have to shape the trajectory, and use it to face the risks.

Such a discourse should be taken into account in setting priorities for investing in AI research and in formulating guidelines, standards and regulatory frameworks, which must be continuously reviewed and updated as we learn more about what is possible, what is necessary and what is to be avoided.

AUTHOR CONTRIBUTIONS

A-MG contributed conception and design of the paper and wrote the manuscript.

ACKNOWLEDGMENTS

The development of this paper benefited from valuable discussions and feedback from my colleagues Roger Bradbury, John Finnigan, Terry Bossomaier, Paul Oppenheimer, and Paulo Santos.

REFERENCES

- Adams, S., Arel, I., Bach, J., Coop, R., Furlan, R., Goertzel, B., et al. (2012). Mapping the landscape of human-level artificial general intelligence. *AI magazine* 33, 25–42. doi: 10.1609/aimag.v33i1.2322
- Burkell, J., and Regan, P. M. (2019). Voter preferences, voter manipulation, voter analytics: policy options for less surveillance and more autonomy. *Internet Policy Review* 8, 1–4. doi: 10.14763/2019.4.1438
- Clark, A. (2015). Radical Predictive Processing. *South. J. Philos.* 53, 3–27. doi: 10.1111/sjp.12120
- Collins, J. (2001). *Good to Great*. New York, NY: Harper Collins.
- Donovan, S. J., Guess, C. D., and Naslund, D. (2015). Improving dynamic decision making through training and self-reflection. *Judgm. Decis. Making* 10, 84–295.
- Dörner, D. (1995). *The Logic of Failure*. New York, NY: Perseus.
- Dörner, D., and Gerdes, J. (2012). “Motivation, emotion, intelligence,” in *International Conference on Systems and Informatics (ICSAI2012)* (Yantai: ICSAI), 691–695.
- Dorner, D., and Güss, C. D. (2013). PSI: a computational architecture of cognition, motivation, and emotion. *Review of General Psychology* 17, 297–317. doi: 10.1037/a0032947
- Dunbar, R. I. M. (1992). Neocortex size as a Constraint on Group Size in primates. *J. Hum. Evol.* 22, 469–493. doi: 10.1016/0047-2484(92)90081-J
- Evans, J., Güss, D., and Boot, W. (2011). Metacognitive prompting aids dynamic decision-making. *Proc. Ann. Meet. Cogn. Sci. Soc.* 33, 3217–3222.
- Grisogono, A. M., and Radenovic, V. (2011). “The Adaptive Stance – steps towards teaching more effective complex decision-making,” in *International Conference on Complex Systems* (Boston, MA: ICCSI).
- Hawkins, J., Lewis, M., Klukas, M., Purdy, S., and Ahmad, S. (2019). A framework for intelligence and cortical function based on grid cells in the neocortex. *Front. Neural Circuits.* 12:121. doi: 10.3389/FNCIR.2018.00121
- Jobin, A., Ienca, M., and Vayena, E. (2019). The global landscape of AI ethics guidelines. *Nat. Mach. Intell.* 1, 389–399.
- Kahneman, D. (2002). *Heuristics and Biases: the Basis of Intuitive Judgment*. New York, NY: Cambridge University Press.
- Kahneman, D., and Tversky, A. (1979). Prospect theory: an analysis of decision under risk. *Econometrica* 47, 263–291.
- Lake, B. M., Ullman, T. D., Tenenbaum, J. B., and Gershman, S. J. (2017). Building machines that learn and think like people. *Behav. Brain Sci.* 40:e253. doi: 10.1017/S0140525X16001837
- Marcus, G. (2009). *Kluge: The Haphazard Evolution of the Human Mind*. New York, NY: Mariner Books.
- Mattingly, H. H., Transtrum, M. K., Abbott, M. C., and Machta, B. B. (2018). Maximizing the Information Learned from Finite data selects a Simple Model. *Proc. Natl. Acad. Sci. U.S.A.* 115, 1760–1765. doi: 10.1073/pnas.1715306115
- Nickerson, R. (1998). Confirmation bias: A ubiquitous phenomenon in many guises. *Rev. General Psychol.* 2, 175–220. doi: 10.1037/1089-2680.2.2.175
- Nowak, M. (2006). *Evolutionary Dynamics - Exploring the Equations of Life*. Belknap Press.
- Rittel, H., and Webber, M. (1973). Dilemmas in a General Theory of Planning. *Policy Sci.* 4, 155–169.
- Staw, B. M. (1997). “The escalation of commitment: an update and appraisal,” in *Organizational Decision Making*, ed Shapira, Zur (New York, NY: Cambridge University Press), 191–215.
- Sterman, J. D. (1989). Misperceptions of feedback in dynamic decision making. *Organ. Behav. Hum. Decis. Process* 43, 301–335. doi: 10.1016/0749-5978(89)90041-1
- Tenenbaum, J. B., Kemp, C., Griffiths, T. L., and Goodman, N. D. (2011). How to grow a mind: statistics, structure, and abstraction. *Science* 331, 1279–1285. doi: 10.1126/science.1192788
- Weick, K. E., and Sutcliffe, K. M. (2001). *Managing the Unexpected*. Jossey-Bass. San Francisco.

Conflict of Interest: The author declares that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2020 Grisogono. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.



Improving the Robustness of Online Social Networks: A Simulation Approach of Network Interventions

Giona Casiraghi and Frank Schweitzer*

Chair of Systems Design, Department of Management, Technology, and Economics, ETH Zurich, Zurich, Switzerland

Online social networks (OSN) are prime examples of socio-technical systems in which individuals interact via a technical platform. OSN are very volatile because users enter and exit and frequently change their interactions. This makes the robustness of such systems difficult to measure and to control. To quantify robustness, we propose a coreness value obtained from the directed interaction network. We study the emergence of large drop-out cascades of users leaving the OSN by means of an agent-based model. For agents, we define a utility function that depends on their relative reputation and their costs for interactions. The decision of agents to leave the OSN depends on this utility. Our aim is to prevent drop-out cascades by influencing specific agents with low utility. We identify strategies to control agents in the core and the periphery of the OSN such that drop-out cascades are significantly reduced, and the robustness of the OSN is increased.

Keywords: socio-technical system, adaptability, robustness, simulations, agent-based model

OPEN ACCESS

Edited by:

Carlos Gershenson,
National Autonomous University of
Mexico, Mexico

Reviewed by:

Radoslaw Michalski,
Wroclaw University of Science and
Technology, Poland
Domenico Rosaci,
Mediterranea University of Reggio
Calabria, Italy

*Correspondence:

Frank Schweitzer
f Schweitzer@ethz.ch

Specialty section:

This article was submitted to
Computational Intelligence in
Robotics,
a section of the journal
Frontiers in Robotics and AI

Received: 01 November 2019

Accepted: 03 April 2020

Published: 28 April 2020

Citation:

Casiraghi G and Schweitzer F (2020)
Improving the Robustness of Online
Social Networks: A Simulation
Approach of Network Interventions.
Front. Robot. AI 7:57.
doi: 10.3389/frobt.2020.00057

1. INTRODUCTION

Self-organization describes a *collective dynamics* resulting from the *local interactions* of a vast number of system elements (Schweitzer, 1997), denoted in the following as *agents*. The macroscopic properties that *emerge* on the system level are often desired, for example, coherent motion in swarms or functionality in gene regulatory networks. But as often these self-organized systemic properties are not desired, for example, traffic jams or mass panics in social systems. Hence, while self-organization can be a very useful dynamics, we need to find ways of controlling it such that systemic malfunction can be excluded, or at least mitigated. This refers to the bigger picture of *systems design* (Schweitzer, 2019): how can we influence systems in a way that optimal states can be achieved and inefficient or undesired states can be avoided?

In general, self-organizing processes can be controlled, or designed, in different ways. On the *macroscopic* or systemic level, *global control parameters*, like boundary conditions, can be adjusted such that phase transitions or regime shifts become impossible. This can be done more easily for physical or chemical systems, where temperature, pressure, chemical concentration, etc. can be fixed. On the *microscopic* or agent level, we have two ways of controlling systems: (i) by influencing agents directly, (ii) by controlling their interactions.

Referring to socio-economic systems, we could, for example, incentivize agents to prefer certain options, this way impacting their utility function. This requires to have access to agents, which is not always guaranteed. For instance, it is difficult to access prominent agents or to influence large multi-national companies. Controlling agents' interactions, on the other hand, basically means to restrict (or to enhance) their communication, i.e., their access to information and dissemination. Restrictions can be implemented both globally and locally.

In this paper, we address one particular instance of social systems, namely *online social networks* (OSN). Prominent examples for such networks are *facebook*, *reddit*, or *Twitter*. OSN are instances of a complex system comprising a large number of interacting agents which represent users of such networks. OSN are, in fact, *socio-technical systems* because they combine elements of a social system, i.e., users communicating, with elements of a technical system, i.e., platforms, protocols, GUI (graphical user interfaces), etc. The technical component is important because it allows to *control* the access to users, as well as their communication. The term *control* refers to the fact that access and interactions are *monitored*, but also *influenced* in different ways.

In reality, it becomes very difficult to control OSN because of their large *volatility*, which has two causes. The first one is the *entry and exit* dynamics, which impacts the number of *agents*: Users enter or leave the OSN at a high frequency. The second one is the *connectivity*, which impacts the number of *interactions*: Users easily connect to and disconnect from other users or interact with lower or higher frequency. They have ample ways of interacting; thus, it becomes very difficult to shield them from certain information.

Because of this volatility, in an OSN *interactions* cannot be fully controlled. But we can certainly influence users via their *utility function*. Users join an OSN for a certain purpose, namely to socialize and to exchange information. Hence, their benefits are a function of the number of other users they interact with. Their costs, on the other hand, result from the effort of maintaining their profile, learning about the features of the graphical user interface, etc. The utility, i.e., the difference between benefits and costs, can then be increased by either increasing the benefits, e.g., by increasing their number of friends, or by decreasing their costs, e.g., by automatizing profile updates, or by a combination of both.

OSN are a paradigm for the emergence of collective dynamics and are much studied because of this. For example, the emergence of trends, fashions, social norms, or opinions occurs as a self-organized process that can sometimes be initiated but hardly be controlled. A worrying trend emerges if users decide to *leave* the social network. If their decision causes other users to leave as well, because they lost their friends, this can quickly result in large drop-out cascades and in the total collapse of the OSN (Kairam et al., 2012). This happened, for example, to *friendster*, an OSN with about 117 million users in 2011. As studied in detail (Garcia et al., 2013), less integrated users left *friendster*, this way, making it less attractive to the remaining users to further stay on the platform.

To model such a self-organized dynamics by means of an agent-based model requires us to solve a number of methodological issues. On the *agent* level, we need to model individual decisions of agents based on their perceived utility, which is to be defined. On the *system* level, we need to quantify how the drop-out of individual agents impact other agents and the whole system, in the end (Jain and Krishna, 1998, 2002). In a volatile system, agents come and go at a large rate, without threatening the stability of the system every time. Hence, we need

to define a macroscopic measure that allows quantifying whether the system is still robust.

Once these methodological issues are solved, we can turn to the more interesting question of *systems design*. This means that, by using our agent-based model, we explore possibilities to influence the system such that it becomes more robust. Our focus will be on the *microscopic* level, i.e., influencing *agents* rather than whole systems. This is sometimes referred to as *mechanism design*. But, different from designing communication, i.e., influencing *interactions*, here we influence agents via their utility functions. This leads to another methodological problem, namely how to identify those agents that are worth to be influenced, i.e., are most promising for reaching a desired system state.

This problem is for networks addressed in the so-called *controllability theory* (Liu et al., 2011), which is very much related to control theory in engineering. It allows to quantify how much of a network is controlled by a given agent, which then can be used to *rank* agents with respect to their *control capacity* (Zhang et al., 2019). To apply this formal framework, however, requires to have a static network, i.e., the interaction *topology* should *not* change on the same time scale as the interaction. So, this framework does *not* allow us to study drop-out cascades in which the network topology changes at every time step. Because of this, in our paper, we have to rely on a *computational approach*, i.e., we use our agent-based model to simulate the decision of agents to leave the network and its impact on the remaining network, while monitoring the overall robustness of the system by means of a macroscopic measure.

With these considerations, we have already specified the structure of this paper. In section 2, we model the decisions of agents and quantify the robustness of the network. In section 3, we introduce a reputation dynamics that runs on the network, to determine the benefits of the agents. In section 4 we highlight the dynamics of the OSN without any interventions, to demonstrate its breakdown. In section 5, eventually, we use our model to explore different agent-based strategies of improving the robustness of the network.

2. ROBUSTNESS OF THE SOCIAL NETWORK

2.1. Agents and Interaction Networks

2.1.1. Networks

For our agent-based model of the OSN we use the specific representation of a *complex network*. The term *complex* refers to the fact that we have a large number of interacting agents such that new system properties can *emerge* as the result of these collective interactions. The term *network* means that agents are represented by *nodes*, and their interactions by *links* of the network. This implies that all interactions are decomposed into dyadic interactions between any two agents.

Using a mathematical language, networks are denoted as *graphs*, nodes as *vertices* and links as *edges*. We can then formally define a *graph* object \mathcal{G} as an ordered pair $\mathcal{G} = \mathcal{G}(V, E)$, where V is the set of vertices of the graph, and E is the set of edges. Vertex

$i \in V$ and $j \in V$ are connected if and only if $ij \in E$. The graph is not static but changes on a time scale T , i.e., $\mathcal{G}(T)$. We call T the *network time* because agents can enter or exit the OSN, this way changing both the number of vertices and edges.

Agents are characterized by a binary state variable $s_i(T) \in \{0, 1\}$, where $s_i(T) = 1$ means that agent i at time T decides to *stay* in the OSN, whereas $s_i(T) = 0$ means that it decides to *leave* the OSN. This decision is governed by a utility function $U_i(T)$:

$$s_i(T) := \Theta[U_i(T)]; \quad U_i(T) = B_i(T) - C_i(T) \quad (1)$$

The Heaviside function $\Theta(x)$ returns 1 if $x \geq 0$ and 0 otherwise. $B_i(T)$ and $C_i(T)$ are the benefits and the costs of agent i at time T . Only if the benefits exceed the costs, agent i will stay in the OSN, otherwise it leaves. The two functions need to be further specified, which is done in section 3.

2.1.2. Interactions

We want to model an OSN; therefore, we consider *directed interactions* between agents. Taking the example of *Twitter*, a directed interaction $i \rightarrow j$ means that agent i is a follower of agent j . Obviously, the reverse does not need to apply but can be frequently observed. Each of these interactions is represented as a directed link in the network \mathcal{G} . A formal expression for the topology of a network with N agents is the *adjacency matrix* $\mathcal{A} \in \mathbb{N}^{N \times N}$ in which the elements a_{ij} are either 0 or 1. This allows to define the *in-degree* d_i^+ and the *out-degree* d_i^- of an agent $i \in V$ as the number of incoming or outgoing links of i . We can also define the *total degree* of agent i as the sum of both in- and out-degree, $d_i = d_i^+ + d_i^-$.

Various works have proposed methods for identifying groups of agents that are stable over time in OSNs. In particular, De Meo et al. (2017) have focused on evaluating the *compactness* of such groups, i.e., the homogeneity in terms of mutual agents' similarity within groups. The concept of compactness, originally introduced in Botafogo et al. (1992), is often used to describe the cohesion of parts of the internet, collaboration networks, and OSNs (Egghe and Rousseau, 2003). Differently from this approach, in this article we aim at characterizing the robustness of the whole network, irrespectively of the stability of specific groups therein. For this reason, we begin our analysis from macroscopic quantities that allow to readily investigate the properties of a complex networks.

The degree distribution is an important macroscopic quantity to characterize a complex network. It is known that OSN have a rather broad degree distribution (Garcia et al., 2013), i.e., many agents are linked to only a few other agents, while a few agents, called hubs, have very many incoming links from other agents. Additionally, OSN often show a so-called *core-periphery structure* (Borgatti and Everett, 2000), in which well-connected agents form a core, whereas agents with only a few, or even no, connections form the periphery. Identifying such structures helps to analyze the robustness of the network. Precisely, we can assume that the OSN is robust, despite an ongoing entry and exit of agents, if the core changes, but continues to exist. This implies that the volatile dynamics mostly affects the periphery. If, however, the drop-out of a few agents is amplified into a large

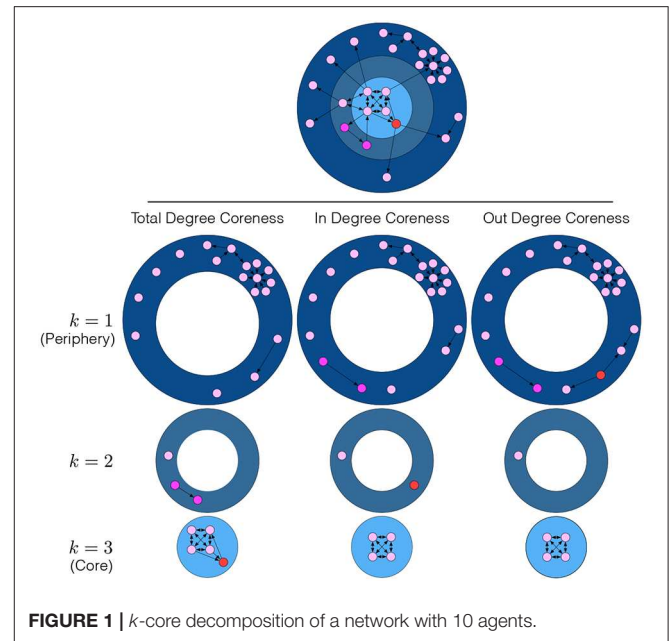


FIGURE 1 | k -core decomposition of a network with 10 agents.

drop-out cascade that affects even the core of the OSN, then the robustness of the system is very low. We need to come up with a robustness measure that reflects such a situation appropriately. This is developed in the next section.

2.2. Quantifying Robustness

2.2.1. Coreness

We decided to use the *coreness* k_i of agents as our starting point because it reflects from a topological perspective how well an agent is integrated into the network (Seidman, 1983). A coreness value k_i allows quantifying the impact on the network when removing agent i . Individual coreness values are obtained by means of a pruning procedure, which is known as *k-core decomposition*. It assigns agents to different concentric shells that reflect the integration of these agents in the network. Specifically, the *k-core* is identified by subsequently pruning all agents with a degree $d_i < k$. Pruning starts with $k = 1$ and stops when all the agents left have a degree greater or equal to k_{\max} . The corresponding *k-shell* then consists of all agents that are in a *k-core* but not in the $(k + 1)$ -core, i.e., agents assigned to a *k-shell* have coreness value $k_i = k$.

Figure 1 provides an illustration of the *k-core decomposition* applied to a network of 10 agents. Agents with a coreness $k_i = 1$ are located in the *periphery* (dark blue), i.e., they are loosely connected with the core. Note that some of these agents have a relatively high degree, in spite of their low coreness. Agents with a coreness $k_i = 2$ are closely connected to, but *not yet* fully integrated into the core, belong to an intermediate shell (blue). The 5 agents with coreness $k_i = k_{\max} = 3$ are the most densely connected ones in this sample network and belong to the innermost core (light blue). This illustrates that the higher the coreness k_i of an agent i , the stronger the impact on the network when removing i because this potentially

disconnects a large number of agents with lower coreness from the network. Conversely, removing agents with low coreness will have a weaker impact on the network because they belong to outer shells, and removing them disconnects a smaller number of agents.

In this article, we want to quantify how much the drop-out of agents will impact the *robustness* of the network. As motivated above, robustness shall be characterized by the *average coreness* of the agents:

$$\langle k \rangle = \frac{1}{N} \sum_{k=1}^{k_{\max}} k n_k ; \quad \sum_{k=1}^{k_{\max}} n_k = N \quad (2)$$

where $N = |V|$ is the total number of (connected and disconnected) agents in the network and n_k is the number of agents with a coreness value $k_i = k$. $\langle k \rangle$ will be high if either most agents have a relatively high coreness, or few agents have a very high coreness. In both cases, the core of the network is less likely to be affected by cascades that started in the periphery. So, $\langle k \rangle$ summarizes the information we are interested in. In this paper, we do not focus on the *heterogeneity* of coreness values, which could be described by the *variance* of the coreness distribution, or by *coreness centralization* (Wasserman and Faust, 1994).

2.2.2. In-Degree and Out-Degree Coreness

The above definition of coreness is based on the total degree d_i of agents, i.e., it is appropriate for *undirected* networks. For the case of a *directed* network discussed in this paper, this may give wrong conclusions about the embeddedness of agents. Therefore, we now introduce two separate measures, in-degree coreness, k_i^+ , and out-degree coreness, k_i^- , which reflect the existence of directed links via the in- and out-degrees d_i^+ , d_i^- .

The results for the different metrics and the differences between them are illustrated in the sample network of 10 agents in **Figure 1**. This network is characterized by 3 k -shells, but it is important to note that the three different coreness metrics possibly assign the same agents to very different k -shells. Take the example of the pair of purple agents that, according to total-degree coreness, are assigned to the shell $k = 2$. If we account for directionality of the links, they are now assigned to $k = 1$, i.e., to the *periphery*. Moreover, the red agent that, according to the total degree coreness, belongs to the core, $k_{\max} = 3$, is now assigned to the shell $k = 2$ if *in-degree coreness* is taken into account, and to $k = 1$, i.e., to the periphery, if *out-degree coreness* is instead considered.

This example makes clear that it very much depends on the *application* whether coreness should be calculated based on directed or undirected links, and whether in- or out-degrees should be considered. In the following we will use *in-degree coreness*, k_i^+ , to compute the average coreness $\langle k \rangle$, Equation (2), i.e., n_k is the number of agents with in-degree coreness $k_i^+ = k$. The reason for this choice comes from the benefits of agents defined in 1 and is discussed in the following section.

3. DYNAMICS ON THE SOCIAL NETWORK

3.1. User Benefits and Costs

To enable a network dynamics on the time scale T , where agents can *leave* the network according to Equation (1), we need to further specify their benefits, $B_i(T)$, and costs, $C_i(T)$. This leads to the question of why, in the real world, users join or leave an OSN. There are certainly different reasons, such as information exchange, maintaining friendship links, or receiving *attention*. From this, we can deduce that benefits should *increase* with the *in-degree* d_i^+ of an agent in a monotonous, but likely non-linear manner. For instance, on Twitter attention increases with the *number* of followers. More important, however, is not just the number, but also the *importance* of the followers. The attention for a user i can considerably increase if it has a number of important users j following. This amplifies the attention because, in an OSN, other users following the important user j this way also receive information from i .

To capture such effects in our agent-based model, we assign to each agent a second state variable, *reputation* R_i , which is continuous and positive. In real-world OSN, user reputation plays an important role and can be proxied by different measures, such as number of likes in Facebook positive votes in Amazon and Dooyoo, or retweets on Twitter. Other proxies take the activity of users into account, for example, the RG score from Researchgate, or the Karma points from Reddit. All of these measures have the drawback that they are (i) specific to the OSN, (ii) depend on the subjective judgment of other users (see e.g., Golbeck and Hendler, 2004, 2006). In the existing literature, the concept of reputation often relates to that of the trust agents pose on each others (Golbeck and Hendler, 2004; Guha et al., 2004; De Meo et al., 2015). Such reputation depends on the activity in the OSN of the agents, e.g., when they evaluate content posted by other agents by “liking” or “disliking” it (Liu et al., 2008; DuBois et al., 2011). In particular, De Meo et al. (2015) have shown that OSN characterized by groups of agents that have higher reputation of each other have higher compactness, and are possibly more stable over time.

Differently from these works, to express agents’ reputation we resort to so-called *feedback centrality* measures. These are prominently known from the early versions of the PageRank algorithm, in which the importance (centrality) of a node in a network entirely depends on the importance of the nodes linked to it. This choice effectively allows us to estimate agents’ reputation directly from the observed topology of the network. This leads to a set of equations for the importance of all nodes that has to be solved in a self-consistent way. While this is a crucial element to define our reputation measure, it is not enough to explain reputation. We also need to consider that reputation fades out over time if it is not continuously *maintained*. Usually, the reputation of an agent can be maintained in different ways, (i) by the own effort of the agent and (ii) by means of direct interactions with others. Such considerations have been formalized in other reputation models (Schweitzer et al., in review). Here, we only consider the increase of reputation coming from other agents, to simplify the formalization.

In the following section, we will specify our dynamics for the reputation of an agent, which leads to a stationary value of $R_i(T)$. Given that we have calculated this value, we posit that the benefit of an agent from being in the OSN comes from its reputation as a good proxy of the attention that this agent receives from others. The absolute value of R_i will also depend on the network size and the density of links. What matters in an OSN is not the *absolute* value, but the reputation of users *relative* to that of others. Therefore, we define the benefit B_i for each agent $i \in V(\mathcal{G})$ as the absolute reputation rescaled by the largest reputation value $R_{\max}(T)$ at the given time T .

$$B_i(T) := b \frac{R_i(T)}{R_{\max}(T)} = b \frac{R_i(T)}{\max_{j \in V(\mathcal{G})} R_j(T)} \quad (3)$$

The constant b allows to weight the benefits from the reputation against the costs.

To specify the costs $C_i(T)$, in our model, we consider two contributions. First, there are fixed costs per time unit, c_0 , that do not depend on the activity of the agents. They capture, in a real OSN, the minimal effort made by users to be present in the OSN, i.e., to learn about the GUI and to maintain the profile. The second contribution comes from the costly interaction with other agents. Because, for instance on *Twitter*, agent i can only control whom to follow, these costs should be proportional to the *out-degree* d_i^- of the agent, $c_i d_i^-$. In a real OSN, the costs per interaction, c_i , are not the same for all users. More prominent users have, for example, much more time constraints because of other activities that compete for their attention. Therefore, it is reasonable to assume that c_i is a non-linear function of the user's reputation, $c_i(R_i) = c_1 R_i^2$. The non-linearity induces a stronger saturation effect for more prominent users in interacting with many other users.

As with the benefits, also the costs should not depend on the absolute reputation of the agent, but on the relative one. This

leads to

$$C_i(T) := c_0 + c_1 d_i^- \left[\frac{R_i(T)}{R_{\max}(T)} \right]^2. \quad (4)$$

Denoting the relative reputation at a given time T as $r_i(T) = R_i(T)/R_{\max}(T)$, we can eventually write down the utility function of agent i , Equation (1), as:

$$U_i(T) = b r_i(T) - c_0 - c_1 d_i^- r_i^2(T) = -c_0 + [b - c_1 d_i^- r_i(T)] r_i(T). \quad (5)$$

3.2. Reputation Dynamics

After linking the utility function of agents to their reputation, we have to specify how to calculate the latter. In accordance with the above discussion, we use the following reputation dynamics:

$$\frac{dR_i(t)}{dt} = -\gamma R_i(t) + \sum_{j \in V(\mathcal{G}(T))} a_{ji} R_j(t) \quad (6)$$

Here, t denotes a *time scale* much shorter than the time scale T at which agents decide whether to stay or to leave the OSN. Hence, compared to the change of the *network*, the change of *reputation* is *fast* enough such that a stationary value $R_i(T)$ is obtained at time T .

The first term in Equation (6) expresses a continuous decay of reputation with a rate γ , to reflect the fact that reputation fades out over time if it is not maintained. The second term captures the increase of reputation coming from other agents linked to agent i , i.e., $a_{ji} = 1$. The summation is over all agents part of the OSN at time T .

Whether or not the reputation values $R_i(T)$ converge to positive stationary values very much depends on the topology of the network expressed by the adjacency matrix \mathcal{A} , as illustrated in **Figure 2**. Specifically, if an agent has no incoming links that boost its reputation, $R_i(t)$ will go to zero. Therefore, even if this agent has an outgoing link to other agents j , it cannot boost

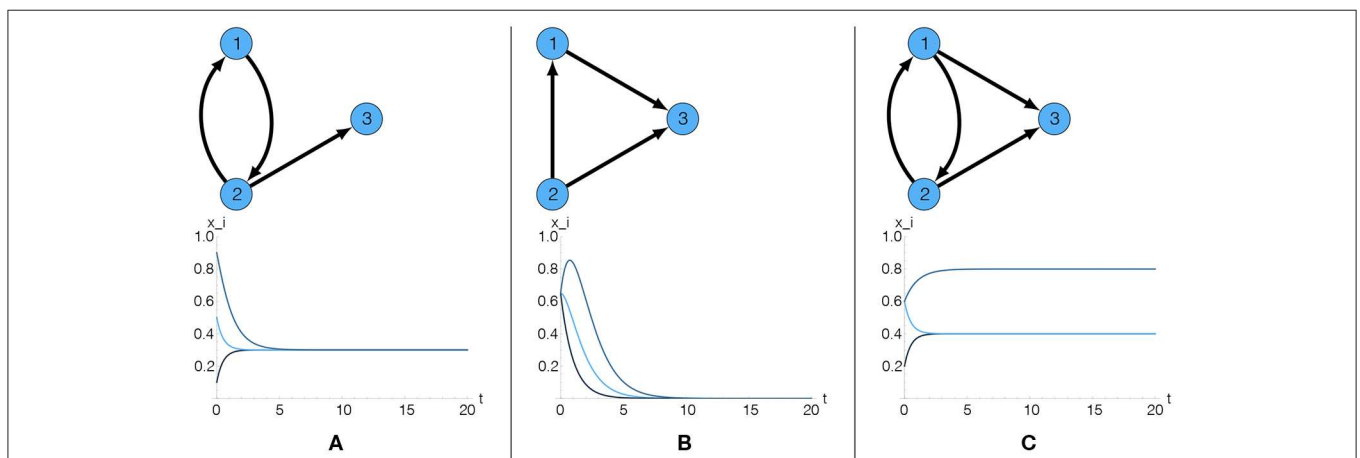


FIGURE 2 | Impact of the adjacency matrix on the reputation $R_i(t)$ of three agents. Only if cycles exist and agents are connected to these cycles, a non-trivial stationary reputation can be obtained. **(A)** The presence of one cycle guarantees a non-trivial stationary reputation, identical for all agents. **(B)** The absence of cycles results in a trivial stationary reputation for all agents. **(C)** The presence of a cycle guarantees non-trivial stationary reputations. Furthermore, two different stationary values appear when agent 3 has 2 incoming links to boost its reputation.

their reputation. Non-trivial solutions depend on the existence of *cycles*, which are formally defined as subgraphs with a closed path from every node in the subgraph back to itself. The shortest possible cycle involves two agents, $1 \rightarrow 2 \rightarrow 1$. This maps to *direct reciprocity*: agent 1 boosts the reputation of agent 2 and vice versa. Cycles of length 3 map to *indirect reciprocity*, for example $1 \rightarrow 2 \rightarrow 3 \rightarrow 1$. In this case, there is *no* direct reciprocity between any two agents, but all of them benefit regarding their reputation because they are part of the cycle. In order to obtain a non-trivial reputation, an agent not necessarily has to be part of a cycle, but it has to be connected to a cycle.

4. DYNAMICS OF THE SOCIAL NETWORK

4.1. Entry and Exit Dynamics

We now have all elements in place to model the entry and exit dynamics of agents in the OSN. At each time step T , agents evaluate their benefits and costs according to Equations (3) and (4). This is based on their relative reputation $r_i(T)$ which has reached a stationary value at time T , according to Equation (6). They then make a (deterministic) decision to either stay or leave the OSN, according to Equation (1).

Hence, at every time T , a number $N^{\text{ex}}(T) < N$ of agents will leave the network. To compensate for this, we assume that the *same* number of new agents will enter the network at the same time, i.e., $N = \text{const.}$ all the time. One may argue that this is at odds with our research question, namely to model how cascades of users leaving impact the robustness of the OSN. But as the empirical case study of the collapse of the OSN Friendster has demonstrated (Garcia et al., 2013), this collapse was *not* due to the fact that no new users entered. Instead, they became *less integrated* into the social network. Signs for this trend became already visible when Friendster had about 80 million users. After that, it still grew up to 113 million users, until it collapsed. So, the problem of the robustness of an OSN cannot be trivially reduced to the (wrong) assumption that there is a lack of new users entering.

Therefore we have to address the question of how, *despite entering of new users*, large drop-out cascades become increasingly likely. To measure the size of the *drop-out cascades*, we will monitor $N^{\text{ex}}(T)$ over time. If this number is consistently large, it becomes evident that even with a large entry rate, new agents cannot substantially stabilize the OSN, hence its robustness is lost. We further need to study how new agents will be *integrated* in the OSN. If at any time T a varying number of $N^{\text{ex}}(T)$ agents *enter*, we have to model how they are linked to the network, to become members of the OSN. We assume that new agents do not have complete knowledge of the network; therefore, to start with, they form *random connections* to a (varying) number of members. Precisely, as in random graphs, new agents create directed links to established agents with a small probability p . Thus, their *expected number* of links is roughly Np .

Because agents leaving delete all their links and agents randomly entering create links, the topology of the network continuously changes at the time scale T . To ensure that the evolution also continues if *no* agent has decided to leave, in this case, we randomly pick one of the agents with the lowest relative

reputation, to replace it with one new agent. To measure how well new agents become integrated into the OSN, we monitor the mean coreness $\langle k \rangle(T)$, Equation (2), over time T . Large values indicate that most agents belong to the core, small values instead that most agents belong to the periphery.

4.2. Results of Computer Simulations

In the following, we discuss the simulation results for a network of fixed size, $N = 20$. Further we use fixed parameters $\gamma = 0.1$, $b = 1$, $c_0 = 0.45$, $c_1 = 0.05$, $p = 0.05$. For a discussion of parameter dependencies and optimal values, see section 5.2.

To initialize our simulations of the network dynamics, we assume that at time $T = 0$, 5 out of 20 agents initially form a fully connected cluster, as shown in **Figure 3A**. This ensures that these five agents have a non-zero reputation at $T = 1$ and thus will not leave the OSN. The remaining 15 agents with reputation zero, however, will be replaced by new agents that randomly create links to the agents in the network. This way, at $T = 50$ already a realistic network structure with a *core*, a *periphery*, different *k-shells* and a few isolated agents emerges, as shown in **Figure 3B**. **Figure 3** displays further snapshots of the network evolution, while the corresponding systemic variables to monitor the dynamics, namely the mean coreness, $\langle k \rangle(T)$, and the number of agents leaving, $N^{\text{ex}}(T)$, are shown in **Figure 4**. From the latter, we can clearly identify three different phases of network evolution.

4.2.1. (I) Build Up Phase

In this initial phase, as already mentioned, the network establishes its characteristic topology. Most agents become tightly integrated into the network, as also visible from **Figures 3B,C**. Because of this, the mean coreness quickly increases, while the number of agents leaving decreases, but both variables show considerable fluctuations.

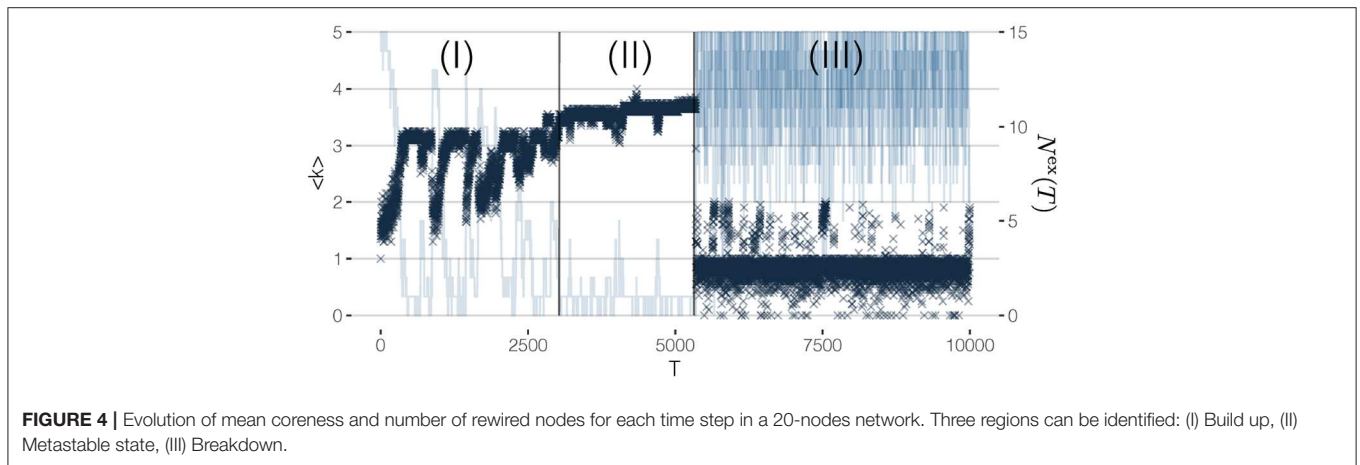
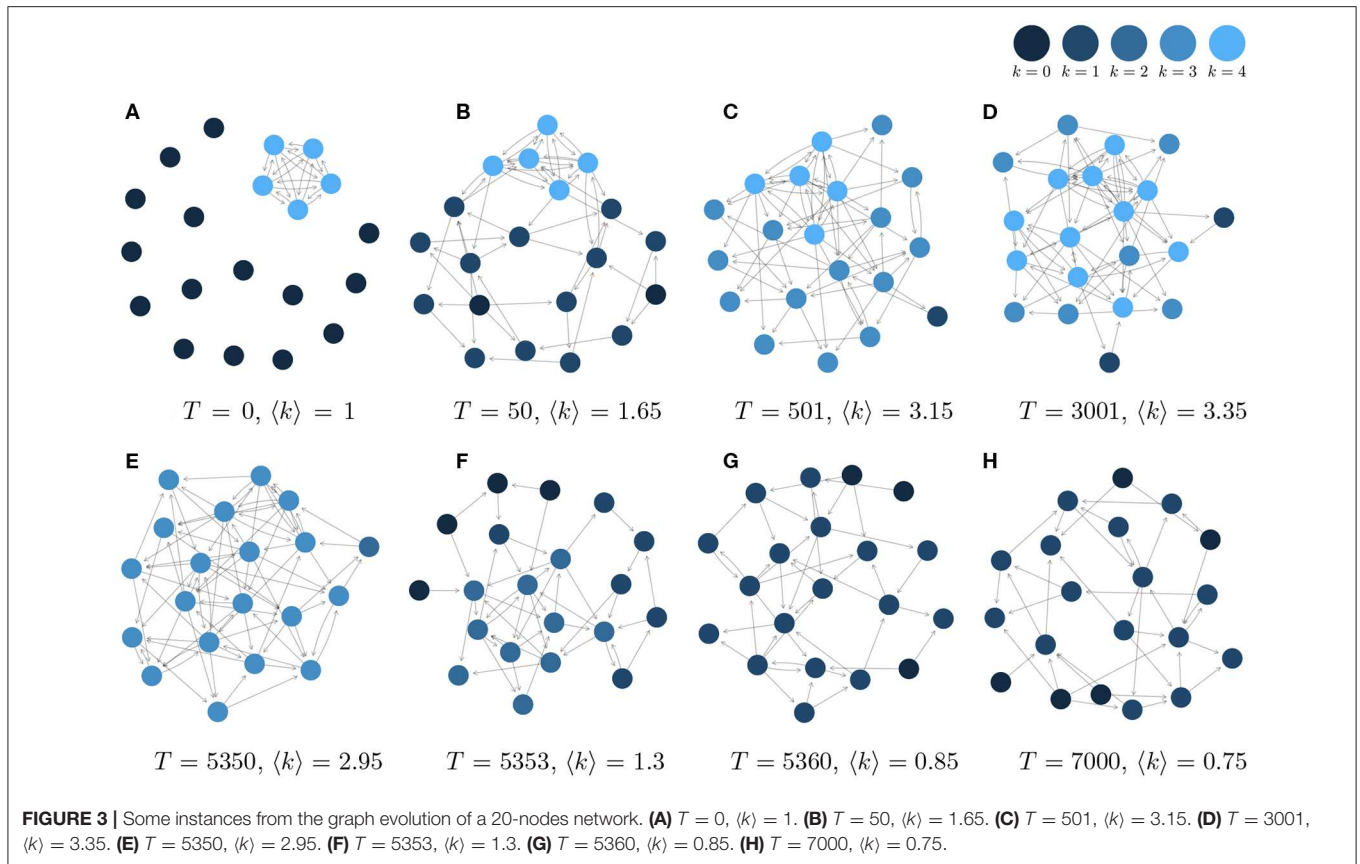
4.2.2. (II) Metastable Phase

After agents have become well-connected to the core, they tend to have higher benefits than costs. If no agent would leave the OSN, we choose one of the agents with the lowest reputation to leave, to keep the network dynamics going. Hence, $N^{\text{ex}}(T) = 1$ or very low, for most of the time, while $\langle k \rangle$ only slightly fluctuates.

Still, the status of the OSN is not stable but only *metastable*, because of the slow dynamics that is illustrated by means of **Figures 3E,F**. Agents that were earlier part of the periphery have now become part of the core, this way *decreasing* the size of the periphery. In fact, the smaller the periphery, the more likely the formation of new links to the core. The probability that a new agent i becomes part of the core Q with size $|Q|$ is given as:

$$P(i \in Q) \geq \binom{|Q|}{k_{\max}^-} p^{k_{\max}^-} \cdot \binom{|Q|}{k_{\max}^+} p^{k_{\max}^+} \quad (7)$$

where k_{\max}^- , k_{\max}^+ are the values for the in-degree and the out-degree coreness of the agents in the *core*. The two r.h.s terms stand for the probability of creating and of receiving links from the core, where p is the probability for an incoming agent to create a new link. $P(i \in Q)$ is indeed increasing with the size of the core, $|Q|$ (Łuczak, 1991).

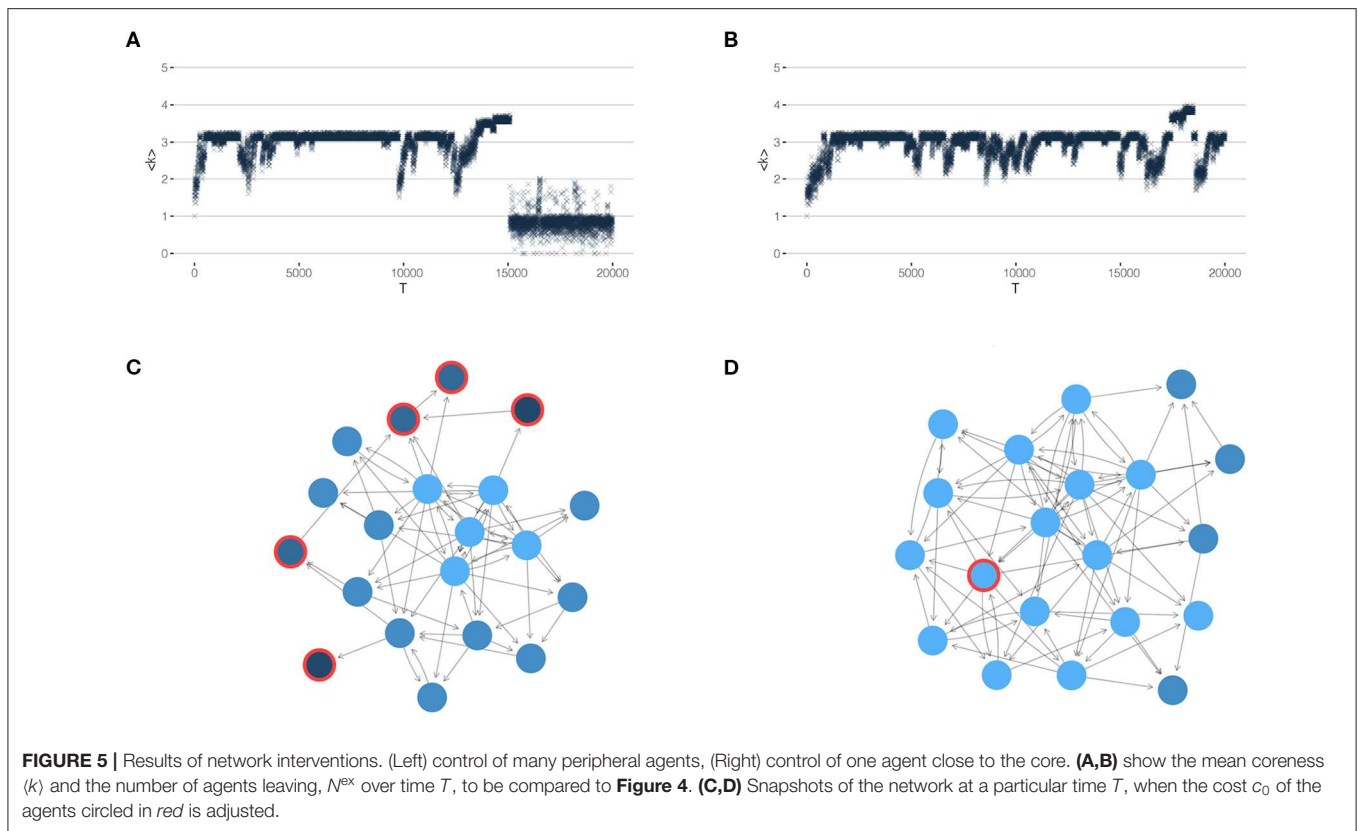


4.2.3. (III) Breakdown Phase

The slow dynamics during phase (II) leads to a point where agents from the outer shells of the in-degree core receive a higher reputation than agents in the core. If no agent decides to leave the OSN, in this situation, an agent from the core is chosen to be removed, because of the lower reputation. This then triggers whole *cascades* of agents leaving, because the drop-out of a core agent abruptly decreases the reputation of other agents in the core and the outer shells. The transition from phase (II) to phase (III) can be seen by the

increasing number of agents leaving, while the mean coreness steadily *decreases*.

Once the core has been destroyed, the OSN has no ability to recover because most agents are replaced at each time step. Nearly all links from the newly entering agents will be to agents from the periphery; thus, the probability of forming a new core is extremely low. The breakdown phase (III) can be characterized not only by the rather low mean coreness and the large number of entries and exits, but also by the much larger fluctuations of both values.



5. IMPROVING ROBUSTNESS

5.1. Network Interventions

The simulation results shown in **Figure 4** make it very clear what we mean by improving robustness: to prevent the *complete breakdown* of the OSN. This does not imply to prevent cascades, which can always happen in response to agents leaving the OSN. But we argue that a social network is *robust* if the decision of agents to leave the OSN will not trigger large cascades of leaving agents that destroy the whole core.

This requires us to influence agents in the OSN such that they decide *not* to leave the network. The trivial solution would be to reduce the costs of *all* agents to a level that always guarantees a positive utility or to increase the benefits in the same manner. A much smarter solution, however, would focus only on a *few* agents, namely those with the ability to prevent large cascades. The problem to identify those agents is addressed in research about network controllability (Liu et al., 2011; Zhang et al., 2016), which is related to control theory. The method assigns a *control signal*, i.e., an incentive to stay or to leave, to the identified agents with the most influence on the network dynamics (Zhang et al., 2019), which are called *driver nodes*. Precisely, this signal is added to the reputation dynamics, Equation (6), of the driver nodes.

We will not follow this formal procedure in our paper for several reasons. The most important one is the continuous evolution of the network topology, which is not considered in the network controllability approach. It would require us to redo the

identification of the driver nodes and the assignment of control signals at every time step T . Further, in our context of users leaving an OSN, these control signals are difficult to interpret because they change the *reputation* dynamics. Our intention instead is to influence the *decisions* of the agents, Equation (1), i.e., to apply control signals to the *costs* of staying in the OSN. Specifically, we apply two different scenarios to incentivize agents (i) from the periphery, or (ii) from the core.

The first scenario is motivated by our insight that large cascades are caused by the *disappearing periphery*. Therefore, a straightforward intervention is to choose agents with a low reputation from the periphery as drivers. These are incentivized to *stay* in the OSN, i.e., their costs are *reduced* such that their utility is increased and they decide *not* to leave. The second scenario is to choose agents close to the core, i.e., from its first outer shells, as drivers. These are incentivized to *leave* the OSN, i.e., their costs are *increased* such that they decide to *not stay*. This more subtle scenario is motivated by the insight that agents that are only close to the core will *not* trigger large cascades if they leave. But if they leave, they considerably reduce the reputation of their closest neighbors, this way *increasing* the size of the periphery. The results of these two scenarios are illustrated in **Figure 5**.

Specifically, in scenario (i), we identify at each time step T all agents from the periphery, i.e., with a coreness value $k_i = 1$. Their cost c_0 is then reduced by 10%, i.e., to $\hat{c}_0 = 0.9c_0$. As **Figure 5A** demonstrates, this scenario can only delay the

complete breakdown (in comparison to **Figure 4** without any interventions). But it cannot completely prevent large drop-out cascade, because the build-up of a large core that eventually gets destroyed is only delayed.

In scenario (ii), on the other hand, we are able to achieve the goal of preventing a complete breakdown. This scenario has remarkable differences to scenario (i): We only incentivize *one* agent, instead of many, and we choose this agent from the vicinity of the *core* instead from the periphery. Precisely, we choose the agent from the first outer shell identified by means of the directed k -core decomposition, i.e., $k_i = k_{\max} - 1$. This agent is enforced to leave by increasing its cost by 10 percent, i.e., to $\hat{c}_0 = 1.1c_0$.

As shown in **Figure 5B**, this scenario considerably improves the robustness of the network, as witnessed by the average coreness. At the same time, because one agent is chosen for control from the beginning, we also observe that the build-up phase (I) is extended in comparison to the case of no control (see **Figure 4**). But phase (II), which was called metastable before, is now considerably extended. We still notice small cascades, but no complete breakdown, i.e., the metastable phase has become a *quasistable* one.

5.2. Life-Time Before Breakdown

The above simulations are both interesting and counter-intuitive because controlling one agent close to the core leads to much better results than controlling many agents from the periphery. We, therefore, continue with a more refined discussion of the peripheral control. As shown, this kind of network intervention increases the time before the breakdown, but cannot completely prevent it. To further quantify this dynamics, we use the *life-time* Ω_Q of the core Q (measured in network time T) as an additional systemic variable (Schweitzer et al., in review). As **Figure 5A** illustrates, for scenario (i) the value of Ω_Q can be clearly obtained from the simulations because of the sharp transition toward the breakdown of the OSN. For scenario (ii), obviously $\Omega_Q \rightarrow \infty$ as **Figure 5B** shows.

We are interested in comparing the life-times of the core for peripheral control and without control (also shown in **Figure 4**). Because Ω_Q changes considerably for different simulations, we use the average life-time $\langle \Omega_Q \rangle$ taken from 100 independent runs with the same setup. We further have to consider that $\langle \Omega_Q \rangle$ depends on other system parameters, notably the system size N . We, therefore, vary N for simulations with peripheral control and without control, keeping all other parameters the same. The results are shown in **Figure 6**, from which we can deduce some interesting insights.

First, we note that for *small* networks ($N < 30$), our peripheral control strategy works very well. The life-times increased considerably in comparison to the no-control reference case. Secondly, we observe that this advantage becomes smaller if the network size increases. For networks larger than $N = 30$, there is almost no difference in life-times between the peripheral control and the no-control case. Further, for $N > 30$ in both cases, the life-time decreases almost linearly with the increasing network size.

The latter observation can be explained from the fact that, with increasing network size N , the network becomes much denser.

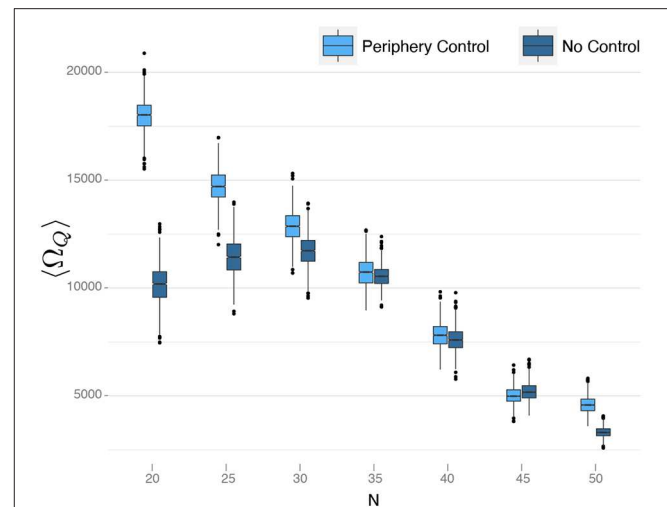


FIGURE 6 | Comparison of different periphery control approaches with fixed control signal. The effectiveness of the control method without adapting the signal to the size of the network decreases with size. In the figure are plotted bootstrap samples for $\langle \Omega_Q \rangle$ obtained from 100 simulations for each network size and each strategy. The control signal used is $u = -0.05$.

We recall that links between agents are formed such that new agents entering the OSN create links to established agents with a *fixed probability*, p . The average number of links per agent is thus Np , i.e., it increases linearly with N . The denser the network, the larger the core and the smaller the periphery. In line with our above discussion, this means less robustness of the network, i.e., the breakdown occurs earlier in time.

The non-monotonous dependence of $\langle \Omega_Q \rangle$ on the network size, for the *no-control* case, results from the fact that the model parameters are not completely independent. This fact is also obvious from Equation (5). Instead, it was already pointed out (Schweitzer et al., under review) that there is an *optimal cost level* to maximize the life-time of the network. This is understandable from our above discussions. If costs are very low, only very few agents will leave the OSN. Because of the slow dynamics described in phase (II), these agents will, at some point, reach a reputation large enough to compare to the core, and hence the core agents will leave. An intermediate cost level, on the other hand, makes sure that this evolution does not take place, or is at least considerably delayed. The optimal cost level that maximizes the life-time, however, also depends on the other parameters, b , N , γ , p .

From **Figure 6**, we can deduce that, for the fixed cost parameters chosen in our simulation, the optimal network size is $N = 30$, simply because, for this size, the life-time is maximized (kept all other parameters the same). Hence, for small networks, $N < 30$, the optimal cost level should be *lower* than what was used in the simulation. Given the suboptimal values, the life-time was also lower for the no-control case. Remarkably, the life-time in case of peripheral control is not affected by this. So, we can conclude that, at least for small networks, peripheral control also compensates for not optimal parameter choices.

For larger networks, $N > 30$, **Figure 6** suggests that there is *no difference* between peripheral control and no control. But this observation is mainly due to the fact that we have not used the optimal parameters for a given network size N . To further investigate this, we have performed an extensive optimization to determine the optimal values for c_0 and \hat{c}_0 for a given N . It then turns out that, with the optimal parameters, the life-times for the peripheral control and no-control cases are no longer the same, but differ *significantly*.

Specifically, we performed two-samples t -tests for the means and *Wilcoxon*-tests for the medians of bootstrap samples of the average life-times (Ω_Q) obtained from the simulations with and without control. As the H_0 hypothesis, we assume that the means of the life-times in both cases are equal and as alternative hypothesis that the life-times are higher in case of peripheral control. Using always the optimal parameters for both cases, we obtained p -values in the order of 10^{-12} for the alternative hypothesis, independent of the network size. This provides strong evidence for the conclusion that the peripheral control *always* improves the robustness of the network, as measured by the life-time before breakdown. For small networks, this holds already for arbitrary parameter choices, for large networks only if the optimal parameters are chosen.

In **Figure 6**, we also plot the bootstrapped 95% confidence intervals for the average life-time (Ω_Q). We note that the size of the confidence interval decreases with N . Hence, for small networks, even optimal parameter values cannot guarantee a minimal variance of Ω_Q , and in single simulations, a breakdown of the network can happen much earlier or later.

Eventually, we also tested whether reputation differences in the peripheral agents matter for the network intervention. While the above simulations assumed that *all* peripheral agents are controlled, we also considered that only peripheral agents with *high*, or with *low* reputation are influenced in their costs. These cases, however, did not generate any remarkable difference with respect to the average life-time.

6. CONCLUSIONS

After more than 35 years of *understanding* complex systems, there should be foundations enough for *managing* them in a better and more quantitative manner. Sadly, to know how systems *work* does not already imply also to know how to *influence* them such that more desired system states are obtained. This holds particularly for socio-economic systems, which are *adaptive*, which means they respond to proposed changes in both intended and unintended ways. *Systems design* (Schweitzer, 2019) therefore has to master a difficult balance: on the one hand, systems should be carefully steered toward a wanted development, on the other hand, systems should not be over-regulated, to not lose their ability to innovate and to find solutions outside the box. This balance cannot be obtained by brute force, in a top-down approach to system dynamics, it has to be found in a bottom-up approach that focuses on the system elements and their interactions.

Our paper contributes to this discussion in several ways. We study a problem of practical relevance that can hardly be solved in a top-down approach: the collapse of an online social network

(OSN) because the decision of some users to leave causes the drop-out of others at large scale. A real-world example is the collapse of the OSN *Friendster* (Garcia et al., 2013). As long as users are free to stay or to leave, the *emergence*, of such large failure cascades cannot be prevented by administrative ruling. Applying global incentives for users to stay, on the other hand, usually implies high costs and questionable efficiency.

Therefore, in this paper, we propose a bottom-up approach to influence the OSN on the level of users, i.e., agents in our model. They can be targeted in two ways: by influencing their interactions or by influencing their utility. We have argued for the latter, because of the large volatility in the dynamics of the OSN. Specifically, we propose to change the costs of particular agents such that the overall robustness of the OSN is increased. As already mentioned in the Introduction, OSN should be seen as *socio-technical systems*, and it is in fact the *technical* component that in principle allows us to influence the costs of users much easier than it would be possible in the offline world.

Improving robustness first requires us to define an appropriate measure of robustness suitable for real-world OSN. Here we propose the *average in-degree coreness*, which does not just reflect the degree of agents but quantifies how well they are *integrated* in the OSN. Next, we have to understand why robustness *decreases* in the absence of network interventions. Based on computer simulations and detailed discussions of agent benefits and costs, we show that it is the changing relation between the core and the periphery of the OSN, which eventually destabilizes the network. Our approach deviates from the one taken in De Meo et al. (2015, 2017) in the fact that we are interested in the robustness of the whole network, and not so much of separate groups. In fact, we learn that is heterogeneity within the network topology, in terms of core-periphery structure, what guarantees robustness. This is in contrast with what expected by generalizing those results obtained for separate groups, where agents' homogeneity increases stability. Moreover, our approach allows to estimate the reputation of agents in the absence of explicit data collecting active declarations of trust between agents in the OSN. To do so, we exploit so-called feedback centralities, that exploit the OSN topology. This is in contrast with common approaches that rely on the presence of likes, dislikes, or agents' ratings to provide a measure for the reputation of agents.

Based on the insights obtained from our analysis, we have proposed two different scenarios for network interventions to improve robustness. The first one targets peripheral agents and reduces their cost, to incentivize them to *stay* in the OSN. The second one targets only *one* agent from a k -shell next to the core and increases its cost, to incentivize it to *leave* the OSN. Both scenarios have in common to increase the size of the periphery, but they reach this goal in different ways. As we demonstrate by means of computer simulations, the first scenario is able to considerably *delay* the breakdown of the OSN, while the second one is able to prevent this breakdown. Dependent on the optimal choice of parameters, we could show that even the peripheral control improves the robustness of the OSN in a *statistically significant* manner. Still, we argue that the second scenario should be the preferred one because it requires (i) to only control a single

agent instead of many, and (ii) less investment because, instead of *decreasing* the costs of many agents via compensations, here the cost is increased.

Our findings are interesting and, at first sight, also counter-intuitive because they challenge our understanding of how to improve the robustness of systems. One could simply argue that the best way to increase robustness is to keep all parts of the system tightly together, to not lose anything. This may apply to mechanical or technical systems. But for socio-technical and socio-economic systems, we have to take into account their adaptivity and their ability to respond to changes in an unintended manner. Therefore, the first step for interventions is to understand the eigendynamics of these systems, i.e., their behavior in the absence of regulations or control. To achieve this understanding in the case of complex systems, agent-based modeling is the most appropriate way. Different from a complex network approach that focuses mainly on the

link topology, agent-based modeling allows also capturing the internal dynamics of the system elements, i.e., the nodes or agents, in response to interactions. Only this advanced level of modeling enables us to propose interventions targeted at specific agents and to investigate how the system as a whole responds to these network interventions.

DATA AVAILABILITY STATEMENT

The raw data supporting the conclusions of this article will be made available by the authors, without undue reservation, to any qualified researcher.

AUTHOR CONTRIBUTIONS

GC and FS designed the research and wrote the manuscript. GC carried out the computer simulations.

REFERENCES

- Borgatti, S. P., and Everett, M. G. (2000). Models of core/periphery structures. *Soc. Netw.* 21, 375–395. doi: 10.1016/S0378-8733(99)00019-2
- Botafogo, R. A., Rivlin, E., and Shneiderman, B. (1992). Structural analysis of hypertexts: identifying hierarchies and useful metrics. *ACM Trans. Inform. Syst.* 10, 142–180. doi: 10.1145/146802.146826
- De Meo, P., Ferrara, E., Rosaci, D., and Sarne, G. M. L. (2015). Trust and compactness in social network groups. *IEEE Trans. Cybernet.* 45, 205–216. doi: 10.1109/TCYB.2014.2323892
- De Meo, P., Messina, F., Rosaci, D., and Sarné, G. M. (2017). Forming time-stable homogeneous groups into online social networks. *Inform. Sci.* 414, 117–132. doi: 10.1016/j.ins.2017.05.048
- DuBois, T., Golbeck, J., and Srinivasan, A. (2011). “Predicting trust and distrust in social networks,” in *2011 IEEE Third Int’l Conference on Privacy, Security, Risk and Trust and 2011 IEEE Third Int’l Conference on Social Computing*, 418–424.
- Egghé, L., and Rousseau, R. (2003). BRS-compactness in networks: theoretical considerations related to cohesion in citation graphs, collaboration networks and the internet. *Math. Comput. Model.* 37, 879–899. doi: 10.1016/S0895-7177(03)00091-8
- Garcia, D., Pavlin, M., and Frank, S. (2013). “Social resilience in online communities: the autopsy of Friendster,” in *Proceedings of the First ACM Conference on Online Social Networks* (New York, NY: Association for Computing Machinery), 39–50. doi: 10.1145/2512938.2512946
- Golbeck, J., and Hendler, J. (2004). “Accuracy of metrics for inferring trust and reputation in semantic web-based social networks,” in *Engineering Knowledge in the Age of the Semantic Web*, eds E. Motta, N. R. Shadbolt, A. Stutt, and N. Gibbins (Berlin; Heidelberg: Springer Berlin Heidelberg), 116–131. doi: 10.1007/978-3-540-30202-5_8
- Golbeck, J., and Hendler, J. (2006). Inferring binary trust relationships in Web-based social networks. *ACM Trans. Internet Technol.* 6, 497–529. doi: 10.1145/1183463.1183470
- Guha, R., Kumar, R., Raghavan, P., and Tomkins, A. (2004). “Propagation of trust and distrust,” in *Proceedings of the 13th Conference on World Wide Web - WWW '04* (New York, NY: ACM Press), 403. doi: 10.1145/988672.988727
- Jain, S., and Krishna, S. (1998). Emergence and growth of complex networks in adaptive systems. *Comput. Phys. Commun.* 122:10. doi: 10.1016/S0010-4655(99)00293-3
- Jain, S., and Krishna, S. (2002). Crashes, recoveries, and core shifts in a model of evolving networks. *Phys. Rev.* 65, 26103–26104. doi: 10.1103/PhysRevE.65.026103
- Kairam, S. R., Wang, D. J., and Leskovec, J. (2012). “The life and death of online groups,” in *Proceedings of the Fifth ACM International Conference on Web Search and Data Mining - WSDM '12* (New York, NY: ACM Press), 673. doi: 10.1145/2124295.2124374
- Liu, H., Lim, E.-P., Lauw, H. W., Le, M.-T., Sun, A., Srivastava, J., et al. (2008). “Predicting trusts among users of online communities,” in *Proceedings of the 9th ACM Conference on Electronic Commerce - EC '08* (New York, NY: ACM Press), 310. doi: 10.1145/1386790.1386838
- Liu, Y.-Y., Slotine, J.-J., and Barabási, A.-L. (2011). Controllability of complex networks. *Nature* 473:167. doi: 10.1038/nature10011
- Łuczak, T. (1991). Size and connectivity of the k-core of a random graph. *Discrete Math.* 91, 61–68. doi: 10.1016/0012-365X(91)90162-U
- Schweitzer, F. (Ed.). (1997). *Self-Organization of Complex Structures: From Individual to Collective Dynamics. Part 1: Evolution of Complexity and Evolutionary Optimization, Part 2: Biological and Ecological Dynamics, Socio-Economic Processes, Urban Structure Formation and Traffic Dynamics*. London: Gordon and Breach.
- Schweitzer, F. (2019). “The bigger picture: complexity meets systems design,” in *Design. Tales of Science and Innovation*, eds G. Folkers and M. Schmid (Zurich: Chronos Verlag), 77–86.
- Seidman, S. B. (1983). Network structure and minimum degree. *Soc. Netw.* 5, 269–287. doi: 10.1016/0378-8733(83)90028-X
- Wasserman, S., and Faust, K. (1994). *Social Network Analysis: Methods and Applications*. Cambridge: Cambridge University Press.
- Zhang, Y., Garas, A., and Schweitzer, F. (2016). Value of peripheral nodes in controlling multilayer scale-free networks. *Phys. Rev. E* 93:012309. doi: 10.1103/PhysRevE.93.012309
- Zhang, Y., Garas, A., and Schweitzer, F. (2019). Control contribution identifies top driver nodes in complex networks. *Adv. Complex Syst.* 22:1950014. doi: 10.1142/S0219525919500140

Conflict of Interest: The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2020 Casiraghi and Schweitzer. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.



Collective Computation in Animal Fission-Fusion Dynamics

Gabriel Ramos-Fernandez^{1,2*}, Sandra E. Smith Aguilar³, David C. Krakauer⁴ and Jessica C. Flack⁴

¹ Departamento de Modelación Matemática de Sistemas Sociales, Instituto de Investigaciones en Matemáticas Aplicadas y en Sistemas, Universidad Nacional Autónoma de México, Ciudad de México, Mexico, ² Unidad Profesional Interdisciplinaria en Ingeniería y Tecnologías Avanzadas, Instituto Politécnico Nacional, Ciudad de México, Mexico, ³ Conservación Biológica y Desarrollo Social A.C., Ciudad de México, Mexico, ⁴ Santa Fe Institute, Santa Fe, NM, United States

OPEN ACCESS

Edited by:

Daniel Polani,
University of Hertfordshire,
United Kingdom

Reviewed by:

Deborah M. Gordon,
Stanford University, United States
Matthew Lutz,
Max Planck Institute of Animal
Behaviour, Germany
Heiko Hamann,
University of Lübeck, Germany

*Correspondence:

Gabriel Ramos-Fernandez
ramosfer@alumni.upenn.edu

Specialty section:

This article was submitted to
Computational Intelligence in
Robotics,
a section of the journal
Frontiers in Robotics and AI

Received: 31 October 2019

Accepted: 05 June 2020

Published: 21 July 2020

Citation:

Ramos-Fernandez G, Smith
Aguilar SE, Krakauer DC and Flack JC
(2020) Collective Computation in
Animal Fission-Fusion Dynamics.
Front. Robot. AI 7:90.
doi: 10.3389/frobt.2020.00090

Recent work suggests that collective computation of social structure can minimize uncertainty about the social and physical environment, facilitating adaptation. We explore these ideas by studying how fission-fusion social structure arises in spider monkey (*Ateles geoffroyi*) groups, exploring whether monkeys use social knowledge to collectively compute subgroup size distributions adaptive for foraging in variable environments. We assess whether individual decisions to stay in or leave subgroups are conditioned on strategies based on the presence or absence of others. We search for this evidence in a time series of subgroup membership. We find that individuals have multiple strategies, suggesting that the social knowledge of different individuals is important. These stay-leave strategies provide microscopic inputs to a stochastic model of collective computation encoded in a family of circuits. Each circuit represents an hypothesis for how collectives combine strategies to make decisions, and how these produce various subgroup size distributions. By running these circuits forward in simulation we generate new subgroup size distributions and measure how well they match food abundance in the environment using transfer entropies. We find that spider monkeys decide to stay or go using information from multiple individuals and that they can collectively compute a distribution of subgroup size that makes efficient use of ephemeral sources of nutrition. We are able to artificially tune circuits with subgroup size distributions that are a better fit to the environment than the observed. This suggests that a combination of measurement error, constraint, and adaptive lag are diminishing the power of collective computation in this system. These results are relevant for a more general understanding of the emergence of ordered states in multi-scale social systems with adaptive properties—both natural and engineered.

Keywords: social systems, distributed computing, inductive game theory, social information, animal foraging, collective intelligence

1. INTRODUCTION

In an influential framework for studying animal social organization, Hinde (1976) stressed that both animal and human societies are multiscale. Short-term interactions between pairs of individuals lead to longer-term social relationships and social structures, with social relationships arising as individuals generalize from a history of social interactions. Hinde noted that individuals

classify social relationships into types (kin, matriline, etc.) regardless of the individuals involved. The idea that primates use abstraction to make sense of their world has been shown in a number of studies subsequent to Hinde (1976) (e.g., Cheney and Seyfarth, 1990, 2008).

Over a series of papers, Flack et al. (Flack, 2012, 2017a,b; Flack et al., 2013; Daniels et al., 2017; Brush et al., 2018) have been developing a theory of collective computation (inspired in part by Hopfield's collective computation in neural networks Hopfield, 1982, 1984; Tank and Hopfield, 1988). In the context of animal behavior, this work links Hinde's (1976) generalization and abstraction processes to the formation of collectives. In Flack and Krakauer's formulation, components (for the purposes of this paper, individuals) reduce uncertainty about the environment or state of a system by coarse-graining fast microscopic behavior (Flack, 2017a). An example of uncertainty reduction would be over the cost of social interaction (Flack, 2012). When coarse-grainings converge (meaning the estimates of regularities are largely shared by individuals), this can produce a coherent mesoscale (e.g., a social network or circuit). This can then function like an information bottleneck (Tishby et al., 2000; Tishby and Zaslavsky, 2015; Flack, 2017a): the strategies, as coarse-grainings, capture regularities individuals perceive in the physical or social environment. The way individuals combine strategies to make decisions in the collective captures the regularities they perceive as most important. Emergent from these slowly changing mesoscopic individual strategies and collective metastrategies is social structure. As a social structure consolidates and individuals start to "reference it" for decision-making, it feeds back through effective downward causation (Flack, 2017a) to modulate the cost of social interaction or interaction with the environment. Once complete, this process can give rise to a new scale, and under suitable conditions, novel functions.

To make this concrete, consider as an example the collective computation of power structure in macaque societies (reviewed in Flack, 2012, 2017a). Individuals summarize fight histories using unidirectional signals. The sender emits the signal once it perceives it is likely to lose a fight. The signal reduces uncertainty in the receiver that the sender agrees to subordination—willingness to yield in future interactions. Encoded in the consolidating network or circuit of signals between group members is information about the distribution of power. Hence the power structure is computed as individuals estimate regularities about fighting abilities and share these opinions with the receiver and other group members via signals. Through this process, different levels of organization arise at successively slower timescales: fights (fast), signaling (slow), and power structure (slowest). The process of generating coarse-grained, slow variables (the signals, properties of the circuits) is the outcome of individual strategic computations (interaction and signaling decisions) that aggregate into an output collectively estimated to fit the state of the environment (Flack, 2017a,b). This two-part process of information accumulation and aggregation makes up collective computation (Daniels et al., 2017; Flack, 2017a).

Among other examples in the animal behavior literature that might result from collective computation are coordinated foraging and predator avoidance in animal groups (Couzin et al., 2003; Gordon, 2016; Sosna et al., 2019), rapid direction changes during collective motion in fish schools and bird flocks (Hein et al., 2015), and distributed foraging in social insects (Gordon, 2016).

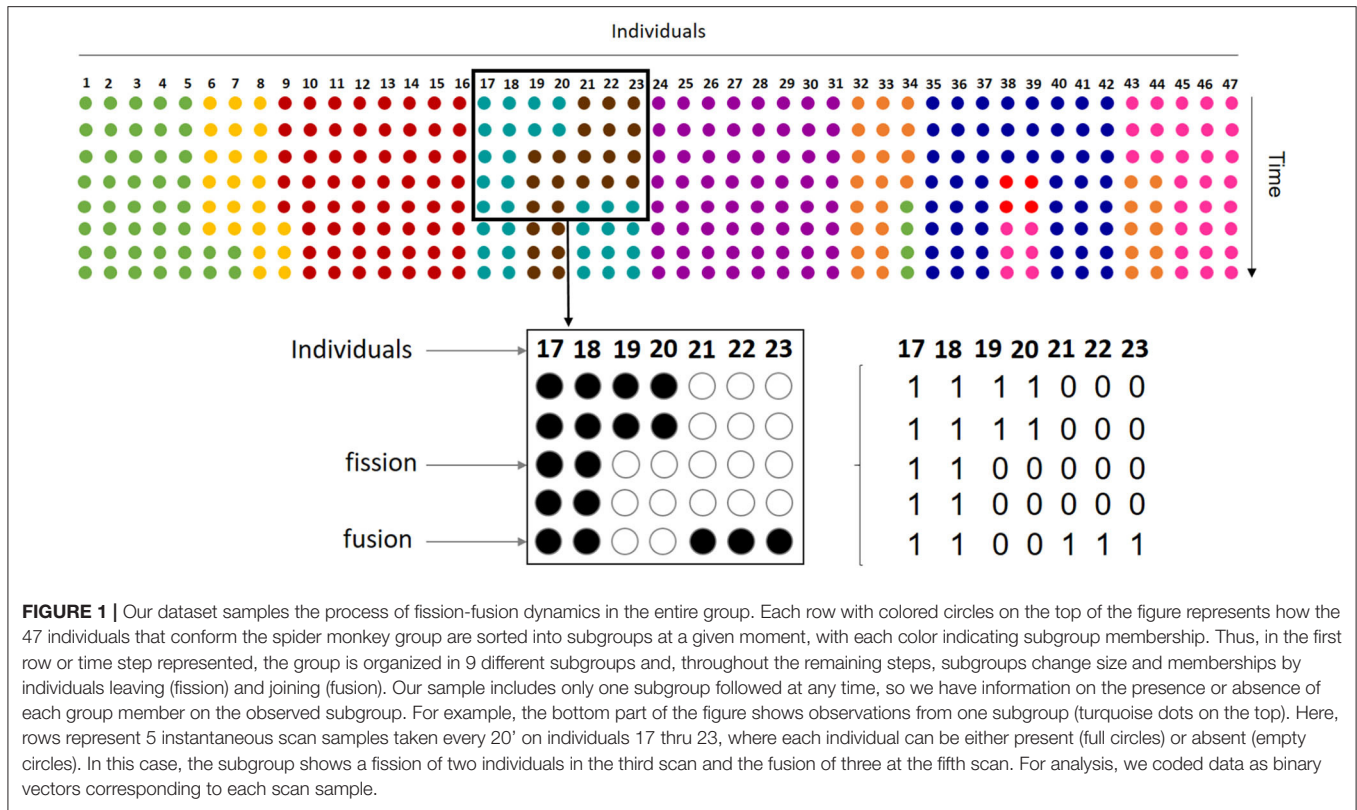
Fission-fusion social dynamics, in which individuals fission and fuse into subgroups of varying size, is a collective pattern arising from individual decisions (Sueur et al., 2011; Ramos-Fernández et al., 2018). These dynamics are thought to be adaptive, as they allow individuals to forage more efficiently in heterogeneous environments, share information about the location of resources, and adjust the size of their subgroups to resource availability (Aureli et al., 2008; Sueur et al., 2011; Palacios-Romo et al., 2019). The individual, strategic decisions to leave or join subgroups, how these decisions influence subgroup size distributions, and whether these are a good fit or even predicted by environmental states, are open questions. Previous work on spider monkeys suggests individuals change their strategies based on environmental states to include the rate at which they encounter fruit and the presence of knowledgeable individuals in social networks (Ramos-Fernández and Morales, 2014; Palacios-Romo et al., 2019).

We study how individual spider monkeys use social knowledge (information accumulation) to collectively compute adaptive subgroup size distributions (information aggregation). We use inductive game theory (DeDeo et al., 2010; Krakauer et al., 2010) to extract stay-leave probabilistic strategies from a time series of subgroup composition. The strategies constitute the microscopic input to the collective computation. From the microscopic input we construct a family of circuits in which nodes correspond to individuals and edges, weighted by probabilities obtained from the data, specify probabilistic rules—strategies—for remaining in or leaving a subgroup. Circuits capture variation in the way individuals integrate over their strategies (see section 3) to decide to stay or go.

Each circuit serves as a mesoscopic hypothesis for how strategies combine to produce decisions and how decisions combine to compute subgroup size distributions. In a computational language, the inputs (individual strategies) combine to produce an output (a subgroup size distribution). We run the circuits forward in simulation to determine how individuals combine strategies and hence how many information sources they take into account to make decisions. We construct a food abundance index based on the size and abundance of fruiting trees and calculate the transfer entropy between this index and the distribution of subgroup size in order to determine whether the circuit that best recovers the observed subgroup size distribution is also optimally computing the state of the environment.

2. DATA

Subgroup composition data were collected in Punta Laguna, Yucatan, Mexico, as part of a long-term study of social behavior



using identified individuals (details about study site and subjects can be found in the **Supplementary Information**). Data consist of scan samples of subgroup composition, taken every 20' during an average of 5 h. per day throughout 2 years (Jan. 2013–Dec. 2014), for a total of 5,780 scan samples. A total of 47 known adult, sub-adult and juvenile individuals were observed during this period (see **Supplementary Table 1**). Thus, each sample is a vector of 47 binary digits, with 0 corresponding to an absence of the individual in the *i*th position and 1 corresponding to a presence (**Figure 1**). Continuous series of scans, averaging 8.4 scan samples (± 3.9 SD), include uninterrupted follows of a subgroup in which at least one individual remained during the full series. Given that the typical duration of a subgroup is 1.5 h. (Pinacho-Guendulain and Ramos-Fernández, 2017), a subgroup may persist over multiple scans. The temporal resolution of this sampling regime was maintained in the analysis in order to obtain a sufficient number of continuous series of scans. Had we resampled the original dataset at a larger temporal scale, we would have lost an important number of continuous series. Also, the persistence of a subgroup over several scans implies that individuals in a subgroup are tolerating one another, which is informative about the weight of their mutual influence (see below).

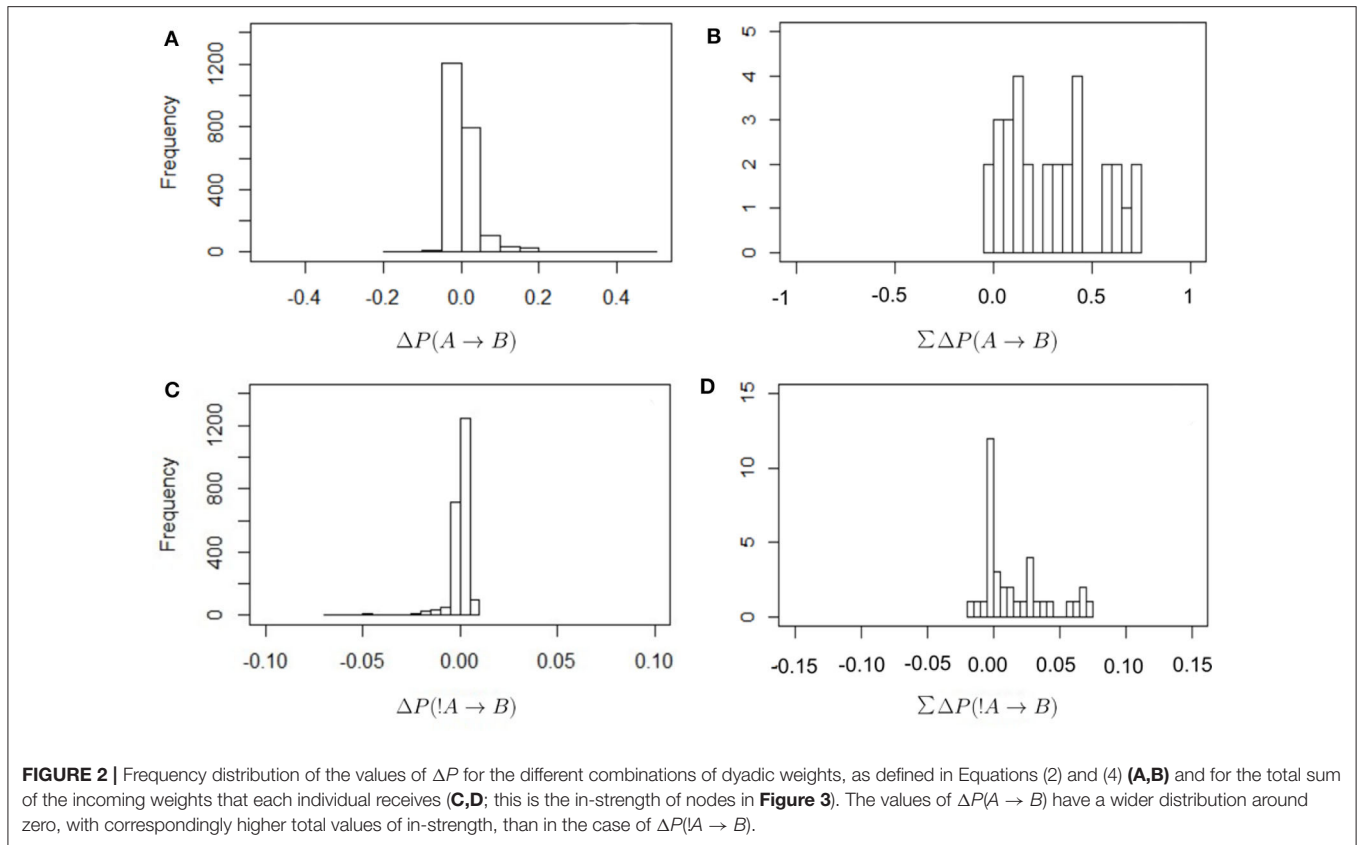
The raw data supporting the conclusions of this manuscript will be made available by the authors, without undue reservation, to any qualified researcher.

3. MICROSCOPIC STRATEGY EXTRACTION AND DISTRIBUTION

We distinguish between strategies and decisions. A decision is binary: to leave or stay in a subgroup (in the original inductive game theory work, to join or avoid a fight, DeDeo et al., 2010). Strategies (called ΔP , as in previous work, DeDeo et al., 2010) are “above-null” probabilities (see below for calculation) describing the weight of individual A’s presence or absence in the *current* subgroup (as determined by scan sampling, see section 2) on individual B’s decision to stay or go from the subgroup in the subsequent sample. Here and in previous work (DeDeo et al., 2010), multiple individuals can influence individual B. Hence B will have multiple strategies and, in the limit, a strategy for every other group member. We address how B integrates strategies to reach a decision in section 4. Here we quantitatively describe how we define and extract strategies from the time series. We end up with a list of pair-wise strategies for which our extraction method indicates above-null support in the time series. We do not consider higher order strategies as in DeDeo et al. (2010).

For all pairs of individuals {A:B, A:C, A:D,...}, we calculate the probability an individual B is present or absent in a sample if individual A was present in the previous sample within the same continuous series of scans:

$$P(A \rightarrow B) = \frac{N(B_{t+1} | A_t)}{N(A)} \tag{1}$$



where $N(B_{t+1} | A_t)$ is the total number of times B was present at time $t+1$ given that A was present at time t within a continuous series of scans and $N(A)$ is the number of times A was present in all samples.

As with previous work (DeDeo et al., 2010), to remove time-independent effects from the transition probabilities (for example, due to general differences in gregariousness), we calculate the difference between the probability inferred from the data and a null expectation:

$$\Delta P(A \rightarrow B) = \frac{N(B_{t+1} | A_t) - N_{null}(B_{t+1} | A_t)}{N(A)}, \quad (2)$$

where $N_{null}(B_{t+1} | A_t)$ is the average number of times B is present at time $t+1$ given that A is present at time t within a continuous series of scans, calculated from 1,000 bootstrapped permutations of the data.

Similarly, we consider the weight of A's absence on the presence of another individual B in a subsequent sample:

$$P(!A \rightarrow B) = \frac{N(B_{t+1} | !A_t)}{N(!A)}, \quad (3)$$

and

$$\Delta P(!A \rightarrow B) = \frac{N(B_{t+1} | !A_t) - N_{null}(B_{t+1} | !A_t)}{N(!A)}, \quad (4)$$

where $N(B_{t+1} | !A_t)$ is the number of times B is present in a sample when A is absent in the previous sample within a continuous series of scans, $N(!A)$ is the number of times A is absent in all samples, and $N_{null}(B_{t+1} | !A_t)$ is the average of the same number for 1,000 bootstrapped versions of the original data.

These ΔP constitute the pair-wise weight of each group member on a given individual's binary decision to leave or join a subgroup.

Figure 2 shows the frequency distribution of the values of ΔP as defined in Equations (2) and (4). In all cases values are centered around zero, with the values of $\Delta P(!A \rightarrow B)$ closer to zero than in other cases. This is because the denominator in Equation (4) is larger than in Equation (2), as it includes all instances of individual A being absent from the observed scan. There are proportionally fewer cases in which B is present after an absence of A because there are many cases where A is absent. Thus, these values of $\Delta P(!A \rightarrow B)$ should be interpreted with care. It is also the case that most values of the total sum of weights received are positive. In other words, most individuals receive a total positive weight from the presence or absence of strategically connected individuals. Only a few cases show a total negative weight of the presence or absence of others.

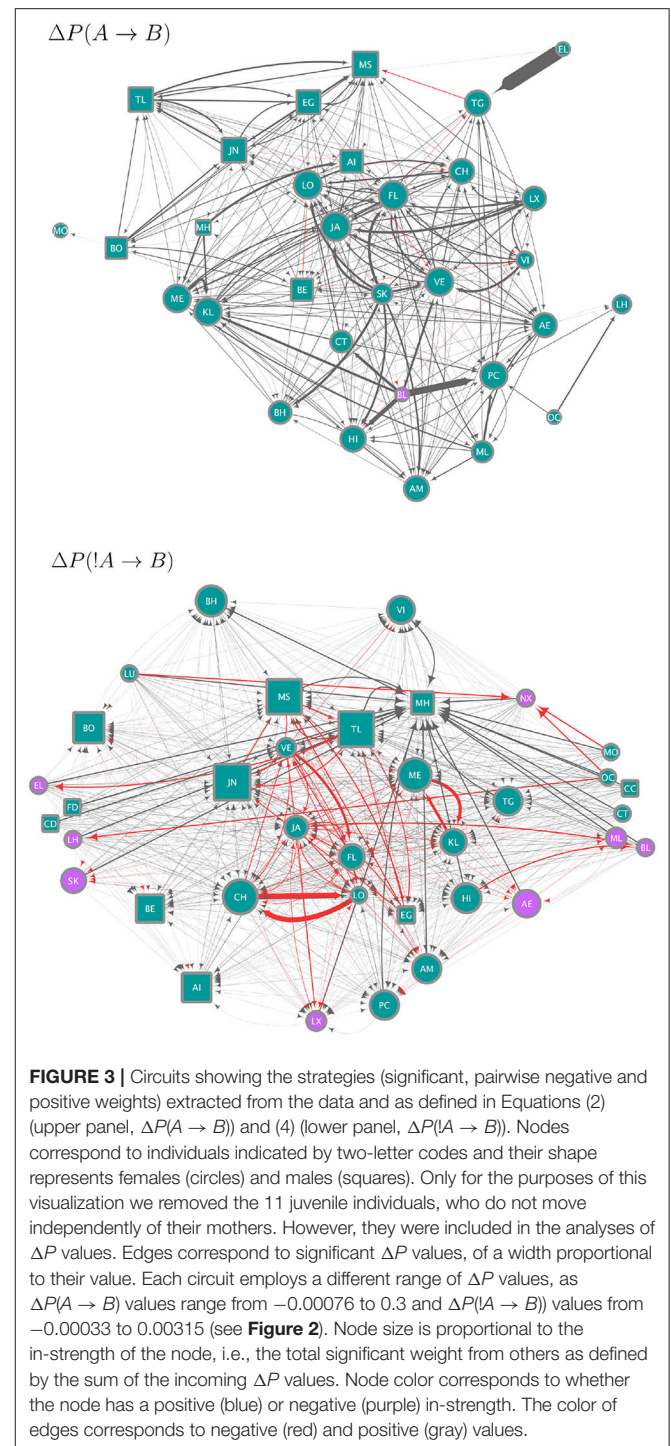
We identified significantly positive dyadic weights as values of ΔP higher than the 95% percentile of the permuted values for each dyad. Accordingly, significantly negative dyadic weights were values of ΔP lower than the 5% percentile of the permuted values for each dyad.

4. MESOSCOPIC CIRCUIT CONSTRUCTION

We use the strategies obtained from the data to construct circuits (*i.e.* the set of all significant ΔP values as weights between all pairs of individuals; this is the mesoscopic level of our analysis) each of which is a hypothesis for (1) how individuals integrate over their strategies to arrive at a binary decision to join or leave a subgroup and, (2) specify how the resulting decisions combine to produce the distribution of subgroup size. The circuits in **Figure 3** give a qualitative summary of significant strategies. For each individual, there are 46 potential weights (significant ΔP values) from either the presence or absence of others at scan time t , which could determine its presence or absence at scan time $t + 1$. The circuits in **Figure 3** show only 31 individual nodes for $\Delta P(A \rightarrow B)$ and 36 for $\Delta P(!A \rightarrow B)$, who were involved in significant weights. On average, each individual in these circuits is linked to 20.25 (± 1.98 SE) other individuals in the $\Delta P(A \rightarrow B)$ and to 31.67 (± 1.40 SE) other individuals in the $\Delta P(!A \rightarrow B)$ circuit (**Figure 3**). Similarly, whereas each of the circuits in **Figure 3** could have up to 1,260 links, the $\Delta P(A \rightarrow B)$ circuit has 314 and the $\Delta P(!A \rightarrow B)$ circuit 570 links. **Supplementary Figure 1** shows the values of all significant weights included in these circuits, as well as the individual in-strength and outstrength.

The circuit for $\Delta P(A \rightarrow B)$ (upper panel in **Figure 3**) represents significant weights of the presence of individual A at scan t on the presence of individual B at scan $t + 1$. Most of the values of $\Delta P(A \rightarrow B)$ were positive or close to zero (see **Figure 2A**), therefore this circuit contains mostly positive weights (gray links), corresponding to weights of attraction. There is an apparent homophily by sex in this circuit, with individuals influencing other individuals of the same sex more than those of the other. Other attractive interactions are those between some pairs of adult females and their subadult daughters (e.g., females VE-VI and JA-LX in the upper panel of **Figure 3**, CH-LO and ME-KL in the lower panel). Individuals differ in their in-strength values (as can be observed in **Figure 2B**) with the individuals with the highest values of in-strength receiving many different weights, some with high values of ΔP , both females and males. Only one individual (female BL) had a negative in-strength value, implying that it received a total negative $\Delta P(A \rightarrow B)$ higher than the total positive $\Delta P(A \rightarrow B)$.

The circuit for $\Delta P(!A \rightarrow B)$ shows a different picture (lower panel in **Figure 3**). Here values were skewed below zero, although overall they were much closer to zero than the values of $\Delta P(A \rightarrow B)$ (**Figure 2**). Even considering that the variation around zero is small, this circuit contains both positive and negative weights, corresponding to repulsion and attraction, respectively, but the most important links are negative or attractive. There is, as in the previous circuit, evidence of some degree of homophily, with individuals of the same sex influencing each other through negative links more than those of the opposite sex. Conversely, a high proportion of positive or repulsive links occur between the sexes. Both males and females have high values of in-strength, although those with a negative in-strength (receiving many negative, attractive weights) in this circuit were all females. Individuals with the highest values of



positive in-strength (corresponding to a total sum of positive or repulsive weights in this network) were males.

Each individual can have multiple strategies, and they can be in conflict (DeDeo et al., 2010), with some weights positive and others negative. In addition, the weight or importance (given by ΔP) of each strategy varies. Hence individuals must integrate over their set of strategies to make a decision about whether

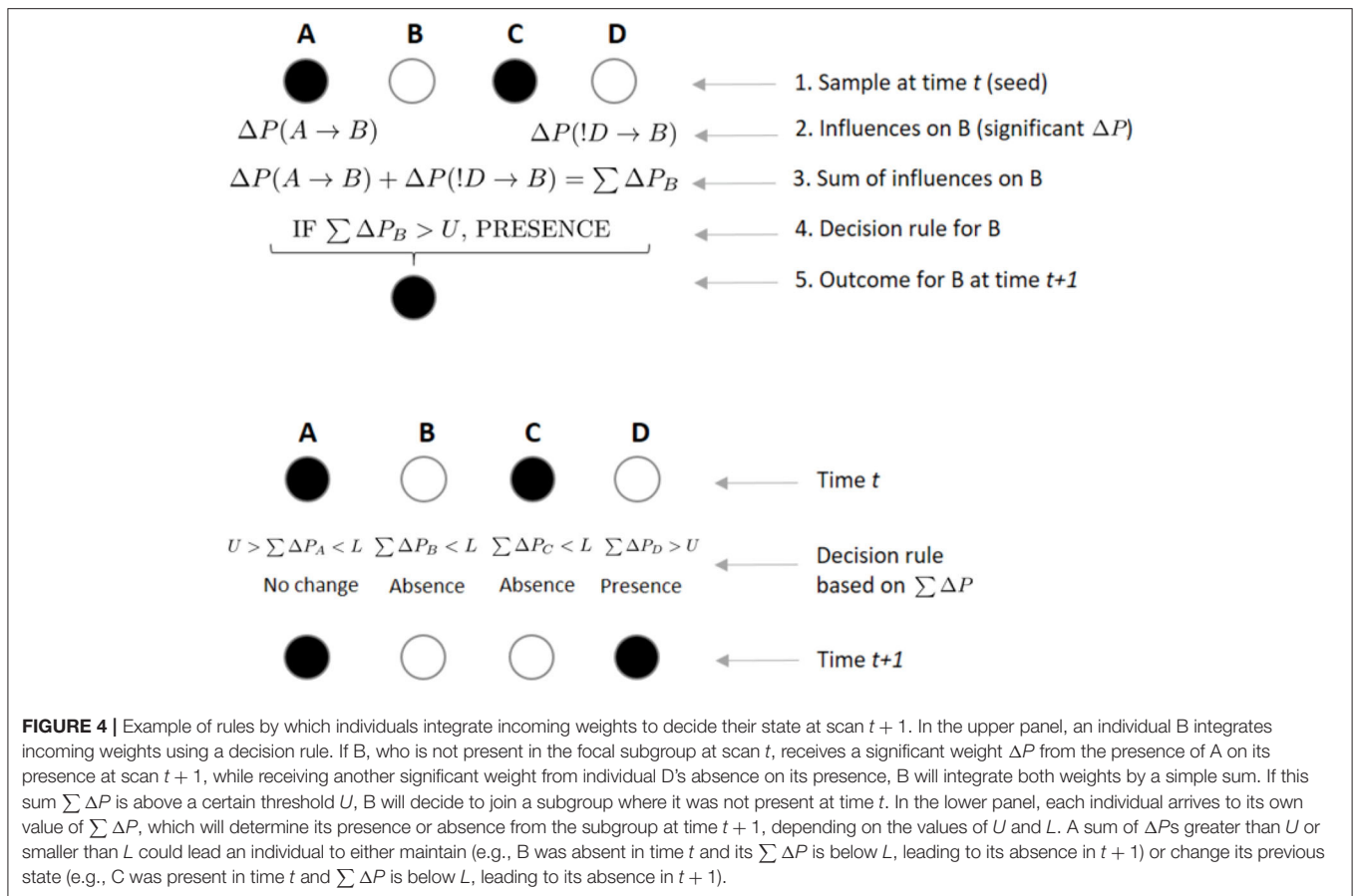


FIGURE 4 | Example of rules by which individuals integrate incoming weights to decide their state at scan $t + 1$. In the upper panel, an individual B integrates incoming weights using a decision rule. If B, who is not present in the focal subgroup at scan t , receives a significant weight ΔP from the presence of A on its presence at scan $t + 1$, while receiving another significant weight from individual D's absence on its presence, B will integrate both weights by a simple sum. If this sum $\sum \Delta P$ is above a certain threshold U , B will decide to join a subgroup where it was not present at time t . In the lower panel, each individual arrives to its own value of $\sum \Delta P$, which will determine its presence or absence from the subgroup at time $t + 1$, depending on the values of U and L . A sum of ΔP s greater than U or smaller than L could lead an individual to either maintain (e.g., B was absent in time t and its $\sum \Delta P$ is below L , leading to its absence in $t + 1$) or change its previous state (e.g., C was present in time t and $\sum \Delta P$ is below L , leading to its absence in $t + 1$).

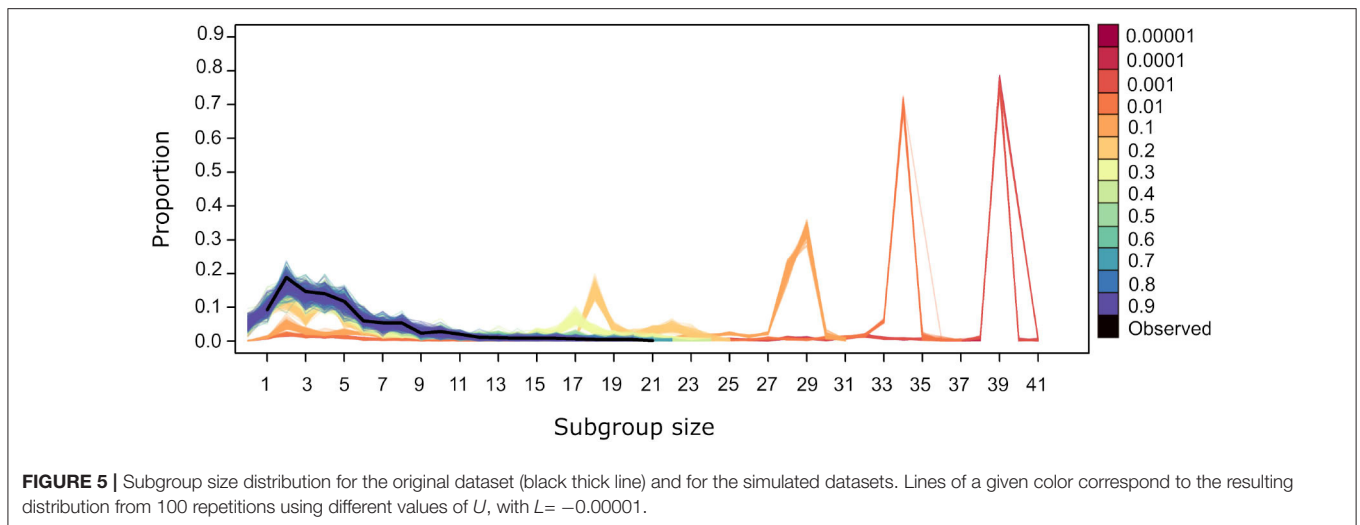
to join or leave the subgroup. **Figures 2B,D** show frequency histograms for these incoming values, corresponding to the in-strength of the nodes in **Figure 3**. These in-strength values can be understood as the likelihood that an individual will be influenced by others: an individual with a high in-strength is more likely to decide to be present due to another individual's presence (in the case of $\Delta P(A \rightarrow B)$ values, upper panel in **Figure 3**) or absence (in the case of $\Delta P(!A \rightarrow B)$ values, lower panel in **Figure 3**) than another individual with a lower in-strength.

We further assume that at any given time t , if the sum of significant ΔP values $\sum \Delta P$ directed toward an individual B is positive and greater than a threshold U , B will be present on the sample at $t + 1$ (irrespective of whether it was present or absent in the previous sample; **Figure 4**). Conversely, if $\sum \Delta P$ is negative and smaller than a threshold L , individual B will be absent from the following sample (again, independently of whether it was present or absent in the previous sample). However, if $L < \sum \Delta P < U$, then there is no effect from others and B remains in the same state as in the previous sample (i.e., present if it was present at time t , absent if it was absent; **Figure 4**). Thus, U is a threshold parameter controlling how likely it is for individuals to be present in a subgroup based on the weight of others. The value L controls the opposite, i.e., how likely it is that individuals will be absent in a subgroup based on the weight of others. Note that the total sum $\sum \Delta P$ includes both the $\Delta P(A \rightarrow B)$

and the $\Delta P(!A \rightarrow B)$ values, such that an individual would be integrating the weights it receives across both circuits shown in **Figure 3**. At higher values of U , the presence of an individual in a subgroup is less likely to be influenced by others. In that sense, high values of U imply less interdependence of individuals in their decisions to be present or not in a subgroup. Conversely, L controls the opposite end of the range of values of $\sum \Delta P$, such that at more negative values of L , an individual should be less likely to be absent from a subgroup due to the previous weight from others. We tested $U = \{0.0001, 0.001, 0.01, 0.1, 0.2, \dots, 0.9\}$ and $L = \{-0.9, -0.8, \dots, -0.1, -0.01, \dots, -0.00001\}$.

Different individuals could actually be using a different value of the U and L thresholds, or the values could change over time, depending on slower ecological variables such as the dry and wet seasons or even longer timescales related to the ecological succession of the forest in the spider monkey's habitat. In this work we assume, as a first approximation, a single value of the threshold parameters for all individuals and seasons.

There are also subtle points here concerning how strategies are aggregated by individuals to produce binary decisions. In previous work (DeDeo et al., 2010), higher order (triadic—C only joins current fight if both A and B were present in the previous fight) as well as pair-wise strategies (A joins if B was previously present) were extracted from time series data and a circuit was constructed for each strategy class. Preliminary



analyses in that work suggest these triadic strategies are non-decomposable into two pair-wise strategies (i.e., not reducible to additive individual or pair-wise interactions; Daniels et al., 2016; Chen et al., 2019). Individuals typically had multiple higher-order strategies and so, as with pair-wise, higher-order strategies were pushed through gates to produce binary decisions. Here we allow for the possibility that individuals take into account multiple strategies and hence be under the influence of multiple individuals, but we do not explore whether the interactions are pair-wise or higher-order.

We use these circuits to generate, by simulation, new datasets from the original dataset. In what follows, we restricted our analyses and simulations to a subset of the original dataset that included the same months for which food abundance data was available (Sep. 2013–Sep. 2014; see section 6), corresponding to 3,032 scan samples. We started by randomly choosing a scan sample (subgroup) that serves as the “seed” or first scan of a sequence of n samples, where n is randomly drawn from the frequency distribution of the number of samples per continuous observation period in the original biweekly period. Thus, the seed establishes which of the 47 monkeys in the group are present or absent in the first sample. Because the seed and the duration of continuous observation periods are selected within observation periods, simulated data contain information about the variation in subgroup size and composition between bi-weekly periods. If an individual A is present in the first scan, the simulation looks at values of $\Delta P(A \rightarrow B)$ and considers any significant values or weights of A on others. If, on the contrary, A is not in the seed, then the simulation looks for significant values of $\Delta P(!A \rightarrow B)$. This applies to all 47 individuals.

These rules are used to determine subgroup composition of the n samples in the continuous observation period. This is repeated for 633 sequences, corresponding to the number of continuous observation periods in the original dataset. In total, we generated 100 simulated datasets for each combination of thresholds U and L .

5. TESTING CIRCUITS IN SIMULATION

Here we assess how individuals integrate strategies to make decisions ΔP and how decisions combine to compute the subgroup size distribution. We do so by asking which circuit, given an integration threshold, produces a simulated data set with a distribution of subgroup size that best recovers the observed one. We used each set of 100 simulated datasets with different values of U to evaluate the set of subgroup size distributions that is in closest correspondence to the observed. We only show the effects of varying U at $L = -0.00001$, since the variation in L for any value of U does not have an effect on the subgroup size distribution. This is likely because values of $\sum \Delta P$ are mostly positive (Figures 2C,D), so very few values are below the L threshold. In other words, even the smallest negative value of L has no effect on the tendency of individuals to modify their presence based on the presence or absence of others.

For values of $U = 0.4$ and above the subgroup size distribution from simulated datasets is similar to the observed (Figure 5). Values of $U < 0.4$ generate distributions where small subgroups are underrepresented and larger subgroups are overrepresented. This is due to the fact that, at lower values of U , individuals are more likely to be influenced by others, both through the significant values of $\Delta P(A \rightarrow B)$ and $\Delta P(!A \rightarrow B)$. The former dominate the dynamics of subgroup size change because they have higher and positive values overall (Figure 2). Thus, when $U < 0.4$, individuals are aggregating more frequently, deciding to join subgroups at higher frequency as in the observed data. Values of $U < 0.4$ give rise to subgroups converging at a single size for each value of U (Figure 5). This may be due to all individuals deciding to join subgroups, even those without significant weights, as must be the case in subgroups larger than 36, the number of nodes in the largest network in Figure 3 that depicts all individuals that are involved in significant weights.

We compared the observed subgroup size distribution and those obtained by simulation under different values of U using

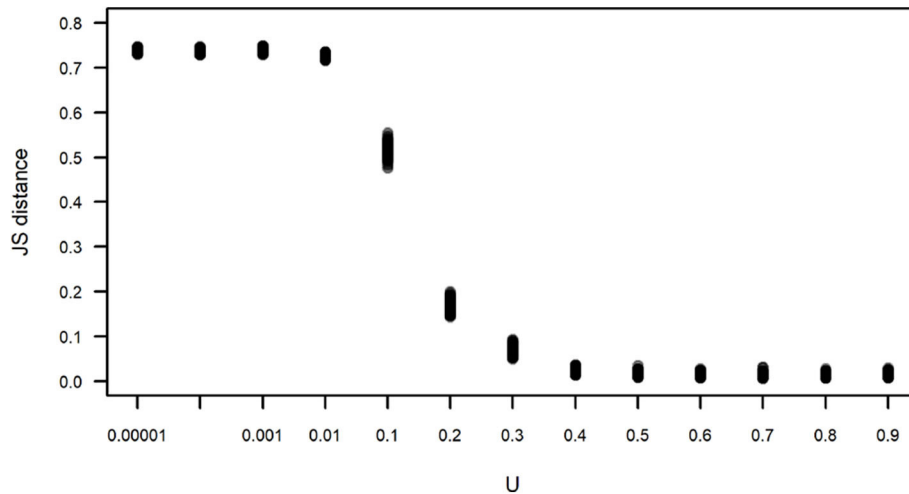


FIGURE 6 | Jensen-Shannon (JS) distance between the observed and simulated subgroup size distributions shown in **Figure 5**. Each dot corresponds to the JS distance between an instance of 100 simulations for each value of U . For all simulations, $L = -0.00001$.

the Jensen-Shannon distance (**Figure 6**). This distance between two random variables x and y is defined as:

$$JS(x|y) = H\left[\frac{x+y}{2}\right] - \frac{1}{2}[H(x) + H(y)] \quad (5)$$

where H is the entropy of each variable, $p(x)\frac{1}{\log p(x)}$ and X and Y are, in this case, the observed subgroup size and the subgroup size obtained in one run of a simulation, respectively. **Figure 6** corroborates what is apparent in **Figure 5**, that simulations run with $U \geq 0.4$ yield subgroup size distributions that are closer and indistinguishable from the observed distribution, with JS values that are close to zero, while simulations run with $U < 0.4$ have an increasing JS with respect to the observed. Simulations run with all values of L for $U=0.4$ yield subgroup size distributions that are equally close to the observed (data not shown).

6. FIT OF OUTPUT TO ENVIRONMENT

A central question is whether the collective computation output is adaptive (Flack, 2017a; Brush et al., 2018). Previous studies of spider monkeys suggest there is a weak relationship between subgroup size and food abundance (Symington, 1988; Pinacho-Guendulain and Ramos-Fernández, 2017). In general, subgroups tend to be larger during periods of high food abundance. This suggests that subgroup size can track the abundance of resources. Here, we investigate whether subgroup size distribution is predicted by the relative abundance of fruiting trees.

We use data from a 1-ha plot where all the trees (diameter at breast height, $D > 10$ cm) from the 15 most consumed species by the monkeys, were monitored bi-weekly for a year from September 2013 to September 2014, comprising 25 monitoring periods. A total of 487 trees were identified, their D was recorded, and every 2 weeks they were assessed for the presence of fruit. The data obtained were used to calculate the proportion of trees with fruit available in a given period expressed in terms of the total tree

D rather than tree number. To do so we calculated the sum of the D values of all the trees with fruit (D_f) in period p divided by the sum of D values for all the trees in the plot (D_i), giving an index of food abundance for a period p , $IFA_p = \sum D_f / \sum D_i$.

Figure 7 shows the time series for the IFA and subgroup size during one year. As mentioned above, maintaining the temporal resolution of the subgroup size time series was important in order to maintain a sufficient number of continuous series of observations. Despite the different temporal resolution of each time series, it seems that subgroup size increases together with IFA during the second wet season.

In previous work, the match between the collective computation output and the environment was evaluated using mutual information (Brush et al., 2018). Here we use transfer entropy:

$$T_{x \rightarrow y}(t) = H(y_t|y_{t-1}) - H(y_t|y_{t-1}, x_{t-1}) \quad (6)$$

This is a measure of how much uncertainty in a variable y is reduced given past states of both y and a variable x that is assumed to be independent of y . This dependence is over and above the uncertainty about y reduced by consideration of its own past state. Here transfer entropy is measuring how much subgroup size uncertainty is reduced by considering past states of subgroup size and IFA, conditioned on the uncertainty reduction by the past states of subgroup size alone. Given the difference in time resolution for the two time series (**Figure 7**), this implies that, within a given bi-weekly period, we are measuring the transfer entropy between a constant value of IFA and varying values of subgroup size. We used the JIDT package (Lizier, 2014) in R (R Core Team, 2017) to estimate the transfer entropy between time series, using the Kraskov estimator with the number of closest neighbors $k = 4$. The two observed time series have a $T_{IFA \rightarrow SGS}(t) = 0.036$ nats.

To explore whether spider monkeys collectively compute a subgroup size distribution that is a good match to the distribution

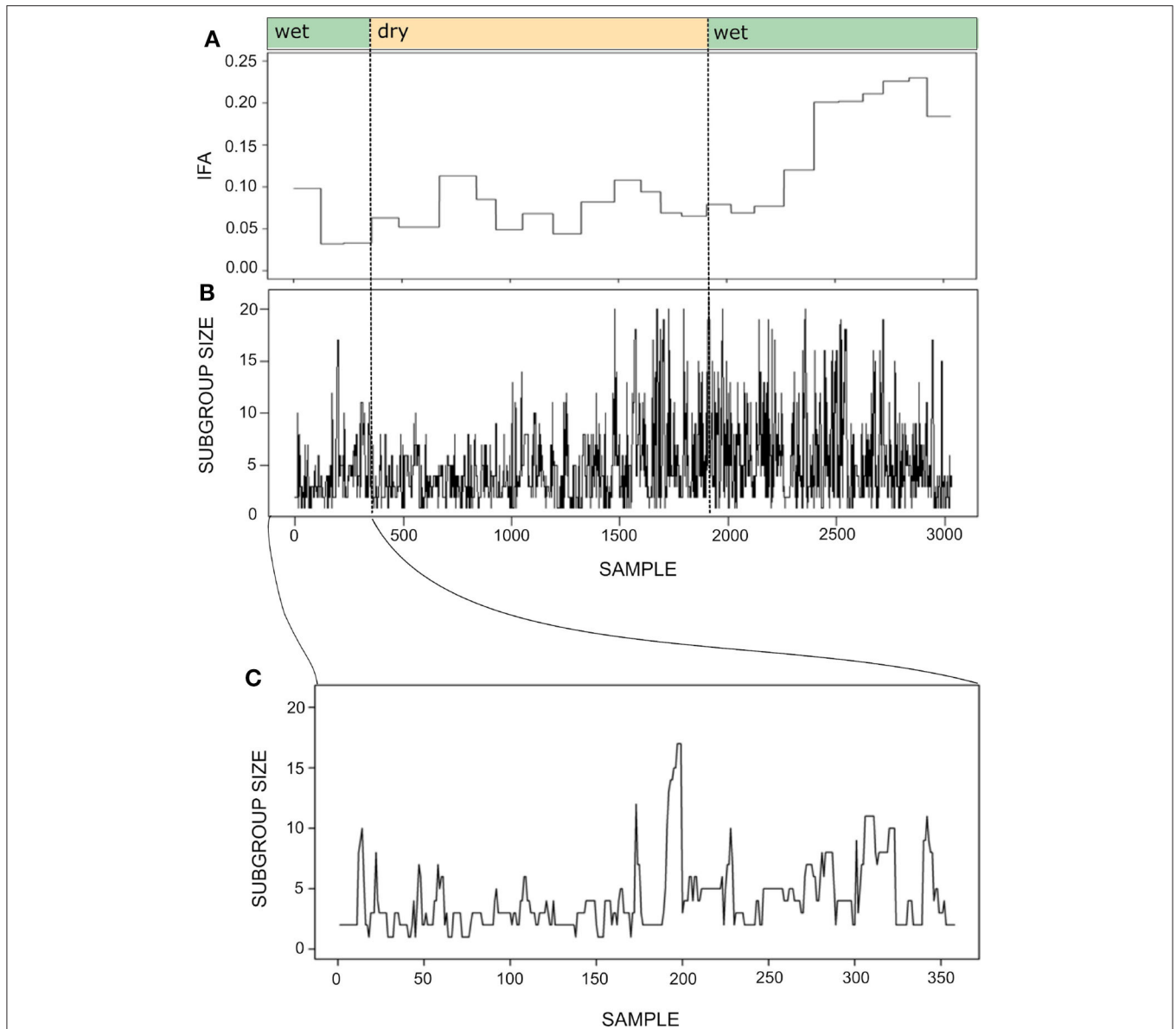


FIGURE 7 | Time series for the index of food abundance (IFA; **A**) and subgroup size (**B**). The IFA measures the overall abundance of fruit in the spider monkey's habitat, considering their most preferred species, their fruiting status and the abundance and relative size of trees (see section 6). The temporal resolution of the subgroup size data is 20 min, whereas food abundance was monitored biweekly. Thus, the IFA series has the same value throughout a given biweekly period, while subgroup size fluctuates at a much finer temporal scale. Noted above are the seasons (wet or dry) to which each sample belongs. Panel (**C**) presents a fragment of the subgroup size time series showing its variation between September 30 and October 31st 2013. Note that the time series was constructed with sets of scan samples taken every 20' collected throughout 4–8 h periods and that subgroups followed in consecutive days were not necessarily the same. Therefore, the spikes and drops observed in the curve do not always reflect fission or fusion events.

of fruiting trees, we assess which of our circuits with different strategy integration rules (described in section 4), computes a distribution of subgroup size that is a good fit to the current abundance of fruiting trees. Shown in **Figure 8** is the time series for the subgroup size values together with the subgroup size time series of all simulated data sets generated for different values of U . **Figure 8** shows what was already apparent in the subgroup size distributions shown in **Figure 5**, but in the

form of a time series: simulated data sets with $U \geq 0.4$ generate a subgroup size distribution that is closest to the observed distribution.

We calculated the transfer entropy between the IFA time series and its corresponding subgroup size time series. We generated simulated data sets that included the same values of IFA as in the original dataset, but because the observation period length could vary (as the length of each observation period, n , was

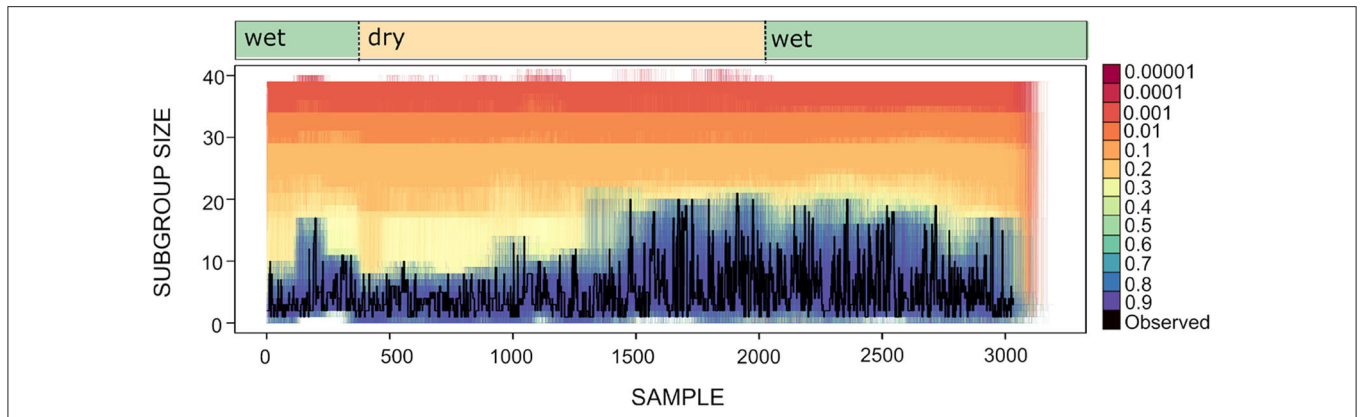


FIGURE 8 | Time series for subgroup size as observed (black line) and simulated (lines of varying color). Each colored line corresponds to an instance of 100 simulations for different values of U and $L = -0.00001$. Wet and dry seasons are noted above.

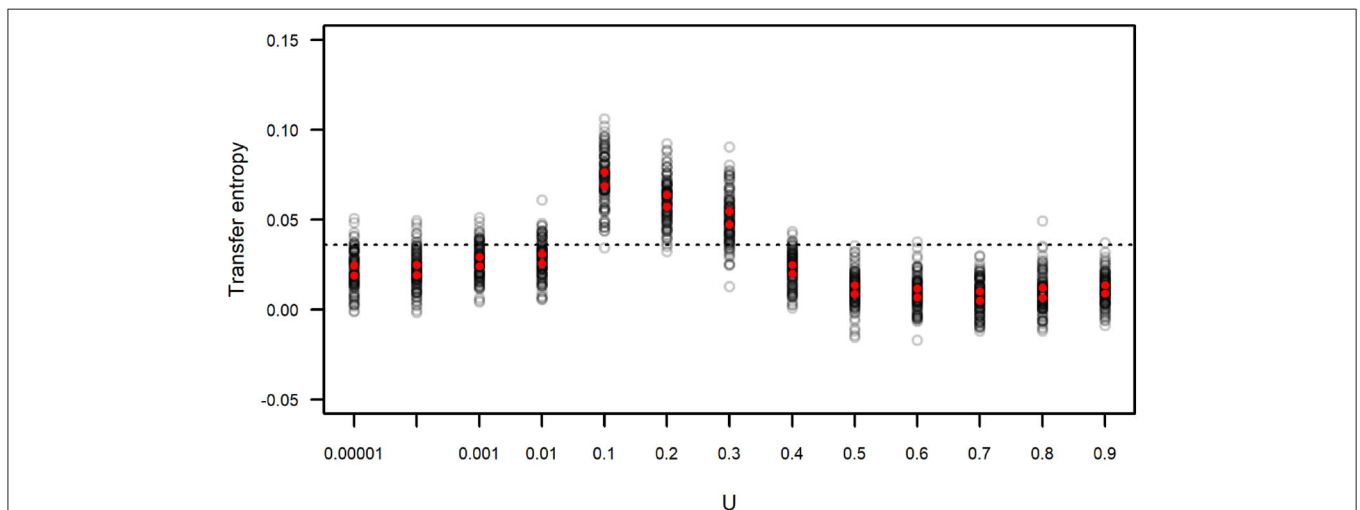


FIGURE 9 | Transfer entropy between simulated IFA and simulated subgroup size. Each gray circle corresponds to an instance of 100 simulations run with varying values of U and $L = -0.00001$. Red dots indicate the upper and lower limits of 99 percent confidence intervals of the mean. The dotted line corresponds to the value of transfer entropy found for the observed IFA and subgroup size data in **Figure 7**.

sampled from the distribution of observed n) there is a certain degree of variation around the observed data. Each simulated IFA series was compared to its corresponding subgroup size series. These values of $T_{IFA \rightarrow SGS}(t)$ are presented in **Figure 9**, which also shows the value of $T_{IFA \rightarrow SGS}(t)$ obtained for the observed IFA and subgroup size time series (**Figure 7**). The results suggest simulated subgroup size data sets with $0.01 < U < 0.4$ match the temporal variation in IFA values better than the empirically observed subgroup size distribution and better than the simulated distributions computed with $U \geq 0.4$.

7. DISCUSSION

Social structure typically changes slowly compared to the interactions giving rise to it. As such, social structure, whether optimal for the environment or not, reduces uncertainty about the future state of the system and provides a relatively stable

background against which individuals can tune their own strategies (Flack, 2017a). Hence there are two challenges for a group computing its social structure: that it changes slowly enough to remain informative for decision-making and that it adaptively tracks the environment.

Frugivorous spider monkeys are faced with two significant sources of uncertainty related to foraging—to discover the location of fruiting trees and to distribute themselves over these fruiting trees to minimize conflict (Aureli et al., 2008) and the costs associated with large groups (Asensio et al., 2009), as well as to maximize resource intake (Symington, 1988). We have used a theory of collective computation (see references in the introduction) to explore how fission-fusion dynamics arises in spider monkey groups and whether the resulting distribution of subgroup size is a good match to the environment. We found spider monkey collectives appear to be able to partially match subgroup size to resource abundance. Our results suggest

however that the collective computation of subgroup size is not optimal with respect to food availability as measured by our index.

In simulating the circuits of subgroup-joining strategies we discover values of a sensitivity parameter U (a measure of the degree of consensus among the incoming weights required for an individual to make a decision about whether to stay or go) leading to a distribution of subgroup size that is a better match (than the observed distribution of subgroup size) to the observed abundance of fruiting trees. This suggests collective computation is under constraint and the system is experiencing adaptive lag—that is, still learning the best collective strategy to integrate information accumulated by group members. The deviation might instead be spurious—an outcome of (1) the way in which we calculate the food abundance index, (2) the fact that the data used to construct the two distributions are noisy and have different time resolutions: food abundance was measured at a bi-weekly scale while subgroup size was observed every 20 min, or (3) other factors besides social knowledge and relationships contributing to subgroup size decision-making.

We should also be cautious in interpreting the power of the collective computation at small U values. In these limits subgroups converge to a constant size where food abundance is expected to be somewhat predictive of size simply because both values remain constant during each bi-weekly period. These caveats aside, whereas collective computation in this system is not optimal, it remains nonetheless predictive and able to capture information about the environment. Specifically, the circuits that capture subgroup joining strategies can aggregate information about the environment. Although we did not study longer timescales, the slowly changing structure of groups provides a means for storing information accumulated by individuals about food availability across years (Palacios-Romo et al., 2019). With individuals that are more than 30 years old (see **Supplementary Information**), who are using spatial memory for their foraging decisions (Valero and Byrne, 2007), the information made available to the group through their experience is likely an important element to track long-term changes in the foraging environment.

Some means by which computations can be refined maximizing the match between group behavior and the abundance of food, includes individuals changing the way they accumulate information and/or compute strategies for staying or leaving, tuning how individuals integrate over those strategies, and tuning how the strategies interact in the circuit to produce subgroup size distributions. For example, are some individuals' strategies (perhaps because they influence many others) exerting a disproportionate effect on the output or do many individuals contribute in small ways? The problem of how collectives achieve optimal information processing is an important one in biology (Tkačik and Bialek, 2016), and near optimal information processing has been discovered in a number of biological systems (e.g., Petkova et al., 2019). However, these examples tend to be relatively simple developmental mechanisms such as segmentation during development of the fruit fly larval body plan. The circuit approach allows the question of tuning to obtain optimal information

processing to be addressed through simulation in more complicated systems.

Additional factors that could affect decision-making, thereby shifting the subgroup distribution from optimal to suboptimal, are a variety of social variables like sex and age, the previous history of interactions, and kinship relationships (Ramos-Fernández et al., 2009; Busia et al., 2017). However, because we are extracting individual strategies directly from the data, these modulating factors are already included in the weights between individuals. Other factors that are currently implicit include the risk of predation or location within the group's home range, which could also affect the subgroup size.

Our results shed light on how a group can best acquire and share information about patchy and dynamic environments. While individual foraging strategies based on spatial knowledge have been well-documented (Janson and Byrne, 2007; Fagan et al., 2013), group foraging strategies are less well-known outside of social insects (Gordon, 2016; cf. Gil et al., 2018). Exchanging information about available patches when foragers disperse and learning about the location and availability of different patches increases the foraging success of the whole group (Falcón-Cortés et al., 2019). The circuit of individual strategies that we infer here is, at least in part, a reflection of information sharing about available patches. Following another individual when ignorant is a simple mechanism of information sharing (Palacios-Romo et al., 2019), that could be reflected in the dyadic weights we have measured. This would lead to a fully connected circuit with information about food sources promoting a flexible grouping pattern that matches heterogeneity in the environment.

It is interesting to compare our approach to that of optimal foraging theory, which would postulate an optimal subgroup size distribution, based on a set of constraints and the best compromise between costs and benefits, which for most cases are unknown (Fretwell and Lucas, 1970; Stephens and Krebs, 1986). An empirical test of this postulate would consist of the match or lack thereof of the observed distribution to the food abundance and this would be interpreted in terms of the unknown mechanisms for how subgroup size comes about (e.g., Chapman et al., 1995). Our approach is more mechanistic: we observe a series of stay-leave decisions resulting from the interactions between individuals and construct a circuit of strategies that serves as a hypothesis for how the subgroup size distribution could emerge. We measure how similar these emerging distributions are to the observed and then test how well the time series matches the environmental variation. That we find alternative circuits that could produce a better match to the environment implies that the system is not necessarily constrained, as would be postulated by optimal foraging theory.

Mutual information, as a measure of uncertainty reduction, has some nice properties. It provides a robust way to study how near optimal a collective behavior is, and this provides a proxy for adaptiveness. We can also study different kinds of uncertainty reduction: an endogenous one, that involves collective computation of social structure that makes the world more predictable for individuals within a system (e.g., Brush et al., 2018); and an exogenous one, whereby collective computation produces social structure that encodes

knowledge about resource availability in the environment (this paper). Uncertainty reduction is consistent with a cost-benefit framework without requiring costs and benefits to be estimated. And quantification of the quality of the output of collective computation in information theoretic terms builds a technical bridge to Boltzmann's and von Neumann's ideas about the role of entropy in generating ordered states (Krakauer et al., 2020) that can form the basis of new levels of individuality, even at the social level.

In addition to assessing whether the output matches the environment, we studied the mechanics of collective computation. Previous work suggests spider monkeys preferentially follow food-aware individuals (Palacios-Romo et al., 2019). In the time series we find evidence in support of this result: we are able to extract significant (above-null) pair-wise probabilistic strategies used by individuals to decide to stay in or leave subgroups. Each individual had 20-30 strategies of varying strength (out of 46 possible). Generally the ΔP were larger for "stay" strategies than "leave" strategies, suggesting possible food presence is a more important factor to spider monkeys than possible food absence. This emphasis on "attraction" might also be important for maintaining cohesion in fission-fusion dynamics in the context of a heterogeneous foraging environment with multiple alternative foraging options (Ramos-Fernández, 2005; Sueur et al., 2011). The strategies we find also recover well-known social patterns for *Ateles* spp., in particular—same sex based homophily for joining and repulsive tendencies between individuals of different sex (Fedigan and Baxter, 1984; Ramos-Fernández et al., 2009). It remains to be determined whether further, more fine-grained patterns like the frequency of dyadic interactions are also recovered by these strategies.

We used the extracted strategies to construct a family of circuits that vary in how individuals integrate these strategies to produce binary decisions to join or leave a subgroup. Individuals can have both repulsion (leave) and attraction (join) strategies. In previous work (DeDeo et al., 2010), strategies were passed through an AND or OR gate that captured conflict averse (all strategies have to say "go" to join a fight) and conflict prone dispositions (one "go" strategy was sufficient to join). Here we use thresholds. To recover the observed subgroup size distribution in simulation requires sums over strategies ($\sum \Delta P \geq U = 0.4$) much larger than the strength of individual strategies (the majority of individual ΔP values are below 0.05). This suggests individual-level decisions, as well as the aggregate output, require that individuals take into account relationships and social knowledge of many group members. If so, this would suggest that spider monkeys rely on social information from the wisdom of crowds (e.g., Jayles et al., 2017; Moreno-Gómez et al., 2017; Kao et al., 2018) to make decisions. These decisions are aggregated to collectively compute subgroup size distributions.

Mesoscale strategic circuits are summaries or average tendencies and therefore provide an economical way to process information. Slow variables, encoded in individual strategies, are compressed summaries of noisy interactions (Flack, 2017b). The idea that the mesoscale circuit is a compressed representation

of microscopic dynamics has parallels in multiplex networks, which have proven to be a better representation of the dynamics of many systems than the simple aggregation of different layers (De Domenico et al., 2015; Smith-Aguilar et al., 2019). Moreover, this way of compressing information may allow the social structure of spider monkeys to be flexible enough to track a dynamic environment, and, at the same time, be robust to disturbances. This has parallels to neural processing (Bassett et al., 2011; Daniels et al., 2017). As we have discussed elsewhere (see Brush et al., 2013, 2018; Daniels et al., 2017; Flack, 2017a) compression and related principles of collective computation have implications for engineered systems, such as web search and swarm robotics (e.g., Bonabeau et al., 1999; Seth, 2001; Young et al., 2013), as well as pattern recognition by artificial neural networks and human reputation networks.

How spider monkeys collectively compute fission-fusion social structure and how these computations can be tuned to realize adaptive variants raises many questions. Using longer time series, we could ask whether collective computation and fit to the environment are being refined and improved over time. With higher resolution data on strategies, and using methods from information theory (e.g., Rosas et al., 2019), it should be possible to quantify the degree to which the output is irreducibly encoded in the circuit as opposed to decomposable. Is social knowledge processed in a pairwise manner or do individuals perceive synergistic interactions among group members (e.g., does individual's A perception of individuals B and C contribute non-additively to its social knowledge)?

Understanding how a natural social system carries out adaptive computations could help to improve the performance of artificial systems. For instance, our results could provide insight into the mechanisms underlying learning through backpropagation in artificial neural networks. The way in which individuals adjust their strategic signaling in computing an appropriate power structure that feeds back to provide information about social interaction cost might be analogous to unsupervised learning (i.e., where the target is endogenous to the system) (Flack, 2017a; Brush et al., 2018). A system like the one we study here, with fission-fusion dynamics that can adjust to environmental conditions like the availability of fruiting trees, might be analogous to supervised learning (i.e., where the target is exogenous to the system). In both cases, feedback might share features with backpropagation in the strong and weak senses—the connection weights in the circuits/networks appear to be adjusted with a combination of vector (Brush et al., 2018) and scalar feedback (Flack et al., 2006) to minimize the network's error function when learning a task (Rumelhart et al., 1986; Lillicrap et al., 2020). This is just one of many exciting comparisons that could be made to better understand how different types of feedback, through tuning (Daniels et al., 2017) and downward causation (Flack, 2017a), shape the ability of the circuit to learn. And, as described in the Introduction, collective coarse-graining can produce a coherent mesoscale functioning as an information bottleneck, an ideal that is at least conceptually similar to the information bottleneck described by Tishby and colleagues to explain how deep neural networks

encode information parsimoniously (Tishby et al., 2000; Tishby and Zaslavsky, 2015; Flack, 2017a).

We have studied how a natural social system collectively computes. This is achieved through feedback among different scales of social organization, as proposed by Hinde's (1976) early paradigm and made explicit in Flack (2017a) and Flack (2017b). Studying collective computation should also find a range of different applications in the engineering of distributed, adaptive systems (Bonabeau et al., 1999).

DATA AVAILABILITY STATEMENT

The datasets generated for this study are available on request to the corresponding author.

ETHICS STATEMENT

The animal study was reviewed and approved by the corresponding authorities in Mexico: the Direccion General de Vida Silvestre, Secretaria de Medio Ambiente y Recursos Naturales.

AUTHOR CONTRIBUTIONS

All authors conceived the idea for this study. JF and DK developed the theory. GR-F and SS designed the data collection and performed the analysis and simulations. SS collected the data. All authors discussed the results and

contributed to the manuscript and giving approval to the final version.

FUNDING

Data collection was aided by a grant from the Mexican 1431 Council of Science and Technology (CONACYT CB157656). JF thanks the Proteus Foundation and the Bengier Family Foundation for support during the project. JF also acknowledges JTF 60501/St. Andrews sub award 13337 for support during the project.

ACKNOWLEDGMENTS

We thank the following people and organizations: Augusto, Macedonio, and Eulogio Canul for their invaluable assistance with data collection; Heiko Hamann, Deborah M. Gordon, and Matthew Lutz for their thorough review of a previous version of this article; Filippo Aureli and Colleen Schaffner for sharing the management of the field project. Instituto Politecnico Nacional and the Center for Complexity Science (C3-UNAM) for their logistical support.

SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/frobt.2020.00090/full#supplementary-material>

REFERENCES

- Asensio, N., Korstjens, A. H., and Aureli, F. (2009). Fissioning minimizes ranging costs in spider monkeys: a multiple-level approach. *Behav. Ecol. Sociobiol.* 63, 649–659. doi: 10.1007/s00265-008-0699-9
- Aureli, F., Schaffner, C. M., Boesch, C., Bearder, S. K., Call, J., Chapman, C. A., et al. (2008). Fission-fusion dynamics: new research frameworks. *Curr. Anthropol.* 49, 627–654. doi: 10.1086/586708
- Bassett, D. S., Wymbs, N. F., Porter, M. A., Mucha, P. J., Carlson, J. M., and Grafton, S. T. (2011). Dynamic reconfiguration of human brain networks during learning. *Proc. Natl. Acad. Sci. U.S.A.* 108, 7641–7646. doi: 10.1073/pnas.1018985108
- Bonabeau, E., Dorigo, M., and Theraulaz, G. (1999). *Swarm Intelligence: From Natural to Artificial Systems*. Oxford: Oxford University Press.
- Brush, E. R., Krakauer, D. C., and Flack, J. C. (2013). A family of algorithms for computing consensus about node state from network data. *PLoS Comput. Biol.* 9:e1003109. doi: 10.1371/journal.pcbi.1003109
- Brush, E. R., Krakauer, D. C., and Flack, J. C. (2018). Conflicts of interest improve collective computation of adaptive social structures. *Sci. Adv.* 4:e1603311. doi: 10.1126/sciadv.1603311
- Busia, L., Schaffner, C. M., and Aureli, F. (2017). Relationship quality affects fission decisions in wild spider monkeys (*Ateles geoffroyi*). *Ethology* 123, 405–411. doi: 10.1111/eth.12609
- Chapman, C. A., Chapman, L. J., and Wrangham, R. W. (1995). Ecological constraints on group size: an analysis of spider monkey and chimpanzee subgroups. *Behav. Ecol. Sociobiol.* 36, 59–70. doi: 10.1007/s002650050125
- Chen, X., Randi, F., Leifer, A. M., and Bialek, W. (2019). Searching for collective behavior in a small brain. *Phys. Rev. E* 99:052418. doi: 10.1103/PhysRevE.99.052418
- Cheney, D. L., and Seyfarth, R. M. (1990). The representation of social relations by monkeys. *Cognition* 37, 167–196. doi: 10.1016/0010-0277(90)90022-C
- Cheney, D. L., and Seyfarth, R. M. (2008). *Baboon Metaphysics: The Evolution of a Social Mind*. Chicago, IL: University of Chicago Press. doi: 10.7208/chicago/9780226102429.001.0001
- Couzín, I. D., Krause, J., et al. (2003). Self-organization and collective behavior in vertebrates. *Adv. Study Behav.* 32, 10–1016. doi: 10.1016/S0065-3454(03)01001-5
- Daniels, B. C., Ellison, C. J., Krakauer, D. C., and Flack, J. C. (2016). Quantifying collectivity. *Curr. Opin. Neurobiol.* 37, 106–113. doi: 10.1016/j.conb.2016.01.012
- Daniels, B. C., Flack, J. C., and Krakauer, D. C. (2017). Dual coding theory explains biphasic collective computation in neural decision-making. *Front. Neurosci.* 11:313. doi: 10.3389/fnins.2017.00313
- De Domenico, M., Lancichinetti, A., Arenas, A., and Rosvall, M. (2015). Identifying modular flows on multilayer networks reveals highly overlapping organization in interconnected systems. *Phys. Rev. X* 5:011027. doi: 10.1103/PhysRevX.5.011027
- DeDeo, S., Krakauer, D. C., and Flack, J. C. (2010). Inductive game theory and the dynamics of animal conflict. *PLoS Comput. Biol.* 6:e1000782. doi: 10.1371/journal.pcbi.1000782
- Fagan, W. F., Lewis, M. A., Auger-Méthé, M., Avgar, T., Benhamou, S., Breed, G., et al. (2013). Spatial memory and animal movement. *Ecol. Lett.* 16, 1316–1329. doi: 10.1111/ele.12165
- Falcón-Cortés, A., Boyer, D., and Ramos-Fernández, G. (2019). Collective learning from individual experiences and information transfer during group foraging. *J. R. Soc. Interface* 16:20180803. doi: 10.1098/rsif.2018.0803

- Fedigan, L. M., and Baxter, M. J. (1984). Sex differences and social organization in free-ranging spider monkeys (*Ateles geoffroyi*). *Primates* 25, 279–294. doi: 10.1007/BF02382267
- Flack, J. C. (2012). Multiple time-scales and the developmental dynamics of social systems. *Philos. Trans. R. Soc. B Biol. Sci.* 367, 1802–1810. doi: 10.1098/rstb.2011.0214
- Flack, J. C. (2017a). Coarse-graining as a downward causation mechanism. *Philos. Trans. R. Soc. A Math. Phys. Eng. Sci.* 375:20160338. doi: 10.1098/rsta.2016.0338
- Flack, J. C. (2017b). “Life’s information hierarchy,” in *From Matter to Life: Information and Causality*, eds S. I. Walker, P. C. Davies, and G. F. Ellis (Cambridge: Cambridge University Press), 283–302. doi: 10.1017/9781316584200.012
- Flack, J. C., Erwin, D., Elliot, T., and Krakauer, D. C. (2013). “Timescales, symmetry, and uncertainty reduction in the origins of hierarchy in biological systems,” in *Evolution of Cooperation and Complexity*, eds K. Sterelny, B. Calcott, and R. Joyce (Boston, MA: MIT Press), 45–74.
- Flack, J. C., Girvan, M., de Waal, F. B. M., and Krakauer, D. C. (2006). Policing stabilizes construction of social niches in primates. *Nature* 439, 426–429. doi: 10.1038/nature04326
- Fretwell, S., and Lucas, H. (1970). On territorial behaviour and other factors influencing habitat distribution in birds. *Acta Biotheoretica* 19, 1–6. doi: 10.1007/BF01601953
- Gil, M. A., Hein, A. M., Spiegel, O., Baskett, M. L., and Sih, A. (2018). Social information links individual behavior to population and community dynamics. *Trends Ecol. Evol.* 33, 535–548. doi: 10.1016/j.tree.2018.04.010
- Gordon, D. M. (2016). The evolution of the algorithms for collective behavior. *Cell Syst.* 3, 514–520. doi: 10.1016/j.cels.2016.10.013
- Hein, A. M., Rosenthal, S. B., Hagstrom, G. I., Berdahl, A., Torney, C. J., and Couzin, I. D. (2015). The evolution of distributed sensing and collective computation in animal populations. *eLife* 4:e10955. doi: 10.7554/eLife.10955
- Hinde, R. A. (1976). Interactions, relationships and social structure. *Man* 11, 1–17. doi: 10.2307/2800384
- Hopfield, J. J. (1982). Neural networks and physical systems with emergent collective computational abilities. *Proc. Natl. Acad. Sci. U.S.A.* 79, 2554–2558. doi: 10.1073/pnas.79.8.2554
- Hopfield, J. J. (1984). Neurons with graded response have collective computational properties like those of two-state neurons. *Proc. Natl. Acad. Sci. U.S.A.* 81, 3088–3092. doi: 10.1073/pnas.81.10.3088
- Janson, C. H., and Byrne, R. (2007). What wild primates know about resources: opening up the black box. *Anim. Cogn.* 10, 357–367. doi: 10.1007/s10071-007-0080-9
- Jayles, B., Kim, H.-R., Escobedo, R., Cezera, S., Blanchet, A., Kameda, T., et al. (2017). How social information can improve estimation accuracy in human groups. *Proc. Natl. Acad. Sci. U.S.A.* 114, 12620–12625. doi: 10.1073/pnas.1703695114
- Kao, A. B., Berdahl, A. M., Hartnett, A. T., Lutz, M. J., Bak-Coleman, J. B., Ioannou, C. C., et al. (2018). Counteracting estimation bias and social influence to improve the wisdom of crowds. *J. R. Soc. Interface* 15:20180130. doi: 10.1098/rsif.2018.0130
- Krakauer, D., Bertschinger, N., Olbrich, E., Flack, J. C., and Ay, N. (2020). The information theory of individuality. *Theory Biosci.* 139, 209–223. doi: 10.1007/s12064-020-00313-7
- Krakauer, D. C., Flack, J. C., Dedeo, S., Farmer, D., and Rockmore, D. (2010). “Intelligent data analysis of intelligent systems,” in *Advances in Intelligent Data Analysis IX*, eds P. R. Cohen, N. M. Adams, and M. R. Berthold (Berlin: Springer), 8–17. doi: 10.1007/978-3-642-13062-5_3
- Lillicrap, T. P., Santoro, A., Marris, L., Akerman, C. J., and Hinton, G. (2020). Backpropagation and the brain. *Nat. Rev. Neurosci.* 21, 335–346. doi: 10.1038/s41583-020-0277-3
- Lizier, J. T. (2014). JIDT: an information-theoretic toolkit for studying the dynamics of complex systems. *Front. Robot. AI* 1:11. doi: 10.3389/frobt.2014.00011
- Moreno-Gómez, S., Sorg, R. A., Domenech, A., Kjos, M., Weissing, F. J., van Doorn, G. S., et al. (2017). Quorum sensing integrates environmental cues, cell density and cell history to control bacterial competence. *Nat. Commun.* 8:854. doi: 10.1038/s41467-017-00903-y
- Palacios-Romo, T., Castellanos, F., and Ramos-Fernandez, G. (2019). Uncovering the decision rules behind collective foraging in spider monkeys. *Anim. Behav.* 149, 121–133. doi: 10.1016/j.anbehav.2019.01.011
- Petkova, M. D., Tkačik, G., Bialek, W., Wieschaus, E. F., and Gregor, T. (2019). Optimal decoding of cellular identities in a genetic network. *Cell* 176, 844–855. doi: 10.1016/j.cell.2019.01.007
- Pinacho-Guendulain, B., and Ramos-Fernández, G. (2017). Influence of fruit availability on the fission-fusion dynamics of spider monkeys (*Ateles geoffroyi*). *Int. J. Primatol.* 38, 466–484. doi: 10.1007/s10764-017-9955-z
- R Core Team (2017). *R: A Language and Environment for Statistical Computing*. Vienna: R Foundation for Statistical Computing. Available online at: <https://www.R-project.org/>
- Ramos-Fernández, G. (2005). Vocal communication in a fission-fusion society: do spider monkeys stay in touch with close associates? *Int. J. Primatol.* 26, 1077–1092. doi: 10.1007/s10764-005-6459-z
- Ramos-Fernández, G., Boyer, D., Aureli, F., and Vick, L. G. (2009). Association networks in spider monkeys (*Ateles geoffroyi*). *Behav. Ecol. Sociobiol.* 63, 999–1013. doi: 10.1007/s00265-009-0719-4
- Ramos-Fernández, G., King, A. J., Beehner, J. C., Bergman, T. J., Crofoot, M. C., Di Fiore, A., et al. (2018). Quantifying uncertainty due to fission-fusion dynamics as a component of social complexity. *Proc. R. Soc. B Biol. Sci.* 285:20180532. doi: 10.1098/rspb.2018.0532
- Ramos-Fernández, G., and Morales, J. M. (2014). Unraveling fission-fusion dynamics: how subgroup properties and dyadic interactions influence individual decisions. *Behav. Ecol. Sociobiol.* 68, 1225–1235. doi: 10.1007/s00265-014-1733-8
- Rosas, F. E., Mediano, P. A. M., Gastpar, M., and Jensen, H. J. (2019). Quantifying high-order interdependencies via multivariate extensions of the mutual information. *Phys. Rev. E* 100:032305. doi: 10.1103/PhysRevE.100.032305
- Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature* 323, 533–536. doi: 10.1038/323533a0
- Seth, A. K. (2001). Modeling group foraging: Individual suboptimality, interference, and a kind of matching. *Adapt. Behav.* 9, 67–89. doi: 10.1177/105971230200900204
- Smith-Aguilar, S. E., Aureli, F., Busia, L., Schaffner, C., and Ramos-Fernández, G. (2019). Using multiplex networks to capture the multidimensional nature of social structure. *Primates* 60, 277–295. doi: 10.1007/s10329-018-0686-3
- Sosna, M. M. G., Twomey, C. R., Bak-Coleman, J., Poel, W., Daniels, B. C., Romanczuk, P., and Couzin, I. D. (2019). Individual and collective encoding of risk in animal groups. *Proc. Natl. Acad. Sci. U.S.A.* 116, 20556–20561. doi: 10.1073/pnas.1905585116
- Stephens, D. W., and Krebs, J. R. (1986). *Foraging Theory*. Princeton, PA: Princeton University Press. doi: 10.1515/9780691206790
- Sueur, C., King, A. J., Conradt, L., Kerth, G., Lusseau, D., Mettke-Hofmann, C., Schaffner, C. M., et al. (2011). Collective decision-making and fission-fusion dynamics: a conceptual framework. *Oikos* 120, 1608–1617. doi: 10.1111/j.1600-0706.2011.19685.x
- Symington, M. M. (1988). Food competition and foraging party size in the black spider monkey (*Ateles Paniscus Chamek*). *Behaviour* 105, 117–132. doi: 10.1163/156853988X00476
- Tank, D., and Hopfield, J. (1988). Collective computation in neuronlike circuits. *Sci. Am.* 257, 104–114. doi: 10.1038/scientificamerican1287-104
- Tishby, N., Pereira, F. C., and Bialek, W. (2000). The information bottleneck method. *arXiv preprint physics/0004057*.
- Tishby, N., and Zaslavsky, N. (2015). “Deep learning and the information bottleneck principle,” in *2015 IEEE Information Theory Workshop (ITW)* (Jerusalem: IEEE), 1–5. doi: 10.1109/ITW.2015.7133169
- Tkačik, G., and Bialek, W. (2016). Information processing in living systems. *Annu. Rev. Condens. Matter Phys.* 7, 89–117. doi: 10.1146/annurev-conmatphys-031214-014803
- Valero, A., and Byrne, R. W. (2007). Spider monkey ranging patterns in mexican subtropical forest: do travel routes reflect planning? *Anim. Cogn.* 10, 305–315. doi: 10.1007/s10071-006-0066-z

Young, G. F., Scardovi, L., Cavagna, A., Giardina, I., and Leonard, N. E. (2013). Starling flock networks manage uncertainty in consensus at low cost. *PLoS Comput. Biol.* 9:e1002894. doi: 10.1371/journal.pcbi.1002894

Conflict of Interest: The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2020 Ramos-Fernandez, Smith Aguilar, Krakauer and Flack. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.



Computational Intelligence for Studying Sustainability Challenges: Tools and Methods for Dealing With Deep Uncertainty and Complexity

Edmundo Molina-Perez^{1*}, Oscar A. Esquivel-Flores² and Hilda Zamora-Maldonado³

¹ Tecnológico de Monterrey, Escuela de Ciencias Sociales y Gobierno, Monterrey, Mexico, ² Institute for Research in Applied Mathematics and Systems, Universidad Nacional Autónoma de México, Mexico City, Mexico, ³ Institute of Economics Research, Universidad Nacional Autónoma de México, Mexico City, Mexico

The study of sustainability challenges requires the consideration of multiple coupled systems that are often complex and deeply uncertain. As a result, traditional analytical methods offer limited insights with respect to how to best address such challenges. By analyzing the case of global climate change mitigation, this paper shows that the combination of high-performance computing, mathematical modeling, and computational intelligence tools, such as optimization and clustering algorithms, leads to richer analytical insights. The paper concludes by proposing an analytical hierarchy of computational tools that can be applied to other sustainability challenges.

Keywords: decision support tools, sustainability, end-of-century climate targets, computational intelligence, climate change, deep uncertainty

OPEN ACCESS

Edited by:

Daniel Polani,
University of Hertfordshire,
United Kingdom

Reviewed by:

Giovanni De Gasperis,
University of L'Aquila, Italy
Holger Lange,
Norwegian Institute of Bioeconomy
Research (NIBIO), Norway

*Correspondence:

Edmundo Molina-Perez
edmundo.molina@tec.mx

Specialty section:

This article was submitted to
Computational Intelligence in
Robotics,
a section of the journal
Frontiers in Robotics and AI

Received: 14 November 2019

Accepted: 16 July 2020

Published: 17 September 2020

Citation:

Molina-Perez E, Esquivel-Flores OA
and Zamora-Maldonado H (2020)
Computational Intelligence for
Studying Sustainability Challenges:
Tools and Methods for Dealing With
Deep Uncertainty and Complexity.
Front. Robot. AI 7:111.
doi: 10.3389/frobt.2020.00111

INTRODUCTION

The resolution of contemporary sustainability challenges requires the consideration of coupled systems, long-term time frames, multiple objectives, and deep uncertainty (Liu et al., 2013, 2016; Hull et al., 2015). For instance, sustainable ecosystem management, water planning, and climate change adaptation and mitigation require the joint consideration of environmental and human systems. These spheres (i.e., systems) are inexorably connected as changes in the behavior and constitution of the natural environment often induce changes in human institutions and incentives. Conversely, the evolution of human preferences, technology, and institutions determines significantly the development trajectories of natural resource systems. Often, if these interactions are not monitored and regulated, one or both systems stop functioning in a sustainable manner (Ostrom, 2009, 2011; Hull et al., 2015). For example, in the context of accelerated global climate change, if anthropogenic emissions continue rising, the growing concentration of greenhouse gases (GHG) in the atmosphere will result in climate imbalances (e.g., changes in precipitation patterns, higher temperatures) that can induce irreversible changes in natural ecosystems (e.g., loss of biodiversity) and in the economy (e.g., higher inequality).

Policy analysis in the context of sustainability is challenging. First, human and environmental spheres are complex systems: path dependencies in both require the consideration of large time frames, and their non-linear interactions induce dynamic behavior that is difficult to anticipate and characterize. Second, deep uncertainty affects both spheres as experts and stakeholders often disagree on the causal representation of these systems, the value of key parameters for analysis, and the relevance of different metrics for describing sustainability (Lempert, 2003; Marchau et al., 2019).

The combination of both conditions, complexity and deep uncertainty, has complicated the role of traditional policy analysis methods when applied to sustainability challenges. On the one hand, the use of simplistic models for analysis can result in omissions relevant for determining long-term outcomes. On the other hand, if the scope of an analysis is too narrow, it is difficult to make the analysis relevant to a wide range of stakeholders (Lempert, 2003; Marchau et al., 2019). Thus, a key emerging question in sustainability sciences is how to design robust policy interventions that explicitly account for complexity and deep uncertainty and which can inform in practical detail public policy discussions of sustainability challenges that affect a wide range of actors.

Modern computational intelligence tools, such as machine learning, optimization, agent-based modeling, and data visualization, offer opportunities for circumventing these limitations (Lempert et al., 2006; Groves and Lempert, 2007; Bryant and Lempert, 2010; Kasprzyk et al., 2013; Isley et al., 2015; Kwakkel, 2017). Yet, their analytical power for sustainability sciences can be best harnessed when these are used in an integrated way. For example, complex simulation models, such as agent-based models (ABMs), can be used as scenario generators in exploratory simulation contexts. Moreover, general purpose and multi-objective optimization techniques can be combined with ABMs to estimate the optimal policy response across large sets of feasible parametrizations. The resulting database can be further analyzed with machine learning algorithms to classify outcomes in terms of the combination of parameter values that trigger different policies. Finally, interactive data visualization techniques can be used to create decision support tools for stakeholders and the public.

Over the last two decades, a growing body of research has applied this integrative approach for studying various sustainability challenges in water (Lempert and Groves, 2010; Groves et al., 2019b; Molina-Perez et al., 2019), energy (Popper et al., 2009), and natural resource planning (Groves et al., 2016; Fischbach et al., 2019). The findings of these studies show that there are no silver bullets for achieving sustainability across human and environmental spheres and that policies that can contribute to achieving sustainable outcomes frequently rely on combinations of different measures that need to be implemented sequentially. First, by addressing immediate vulnerabilities through robust policies. Second, by responding adaptively to medium and long-term changes in both spheres (Groves et al., 2019b; Molina-Perez et al., 2019). This body of research, defined as Decision Making under Deep Uncertainty (DMDU) (Marchau et al., 2019), has cemented the foundations for the general application of computational intelligence tools to sustainability sciences.

This paper applies DMDU methods—specifically Robust Decision Making (Lempert, 2003; Groves et al., 2019b)—to structure an analysis of global climate change mitigation and to demonstrate that the combination of multiple computational

tools for analyzing this sort of sustainability challenges leads to richer analytical insights than those produced by traditional monodisciplinary studies. Particularly, our analysis shows that by integrating optimization, integrated assessment models, and machine learning algorithms, it is possible to quantitatively identify key drivers of vulnerability of climate change mitigation policies. It also shows that alternative policy proposals can work as complements across regions to cost-effectively decarbonize the global economy. The paper concludes by proposing an analytical hierarchy of computational tools that can be applied to other sustainability challenges.

COMPUTER MODELING FOR CLIMATE CHANGE POLICY ANALYSIS

Virtual Laboratories and Policy Regimes

Simulation models are popular tools in the field of climate change because of (a) the long-term time horizons needed to be taken into consideration, (b) the heterogeneous economic and technological conditions of countries and industries, and (c) the non-linearities and path dependencies associated with climate policy. To highlight how the combined use of integrated assessment models (IAMs) and other computational intelligence tools can result into a more detailed understanding of sustainability challenges, in this study we use the Exploratory Dynamic Integrated Assessment Model (EDIAM) developed by Molina-Perez (2016).

The EDIAM model is primarily based on the theoretical framework developed by Acemoglu et al. (2012), which takes into account the interrelation between climate mitigation, innovation, and growth. Particularly, it describes the propagation of climate policy impacts in the economy through endogenous productivity changes that affect labor, energy, and technology markets. EDIAM expands over this framework by including the role of learning-by-doing, differentiating technology properties across sectors, modeling entrepreneurs' investment decisions in continuous form, considering the role of technological transferability across nations, and calibrating environmental equations using a full ensemble of Coupled Model Intercomparison Project Phase 5 (CMIP5) climate projections (Taylor et al., 2012; IPCC, 2013).

The motivation for using EDIAM as an instrument for experimentation in this paper is based on four of its characteristics. First, it emphasizes the role that technology policy plays in climate mitigation. Second, it describes how climate policy propagates through time, changing the incentives of economic agents (i.e., path dependency). Third, it considers the interconnection between regions and between the environment and the economy and its role in shaping global outcomes (i.e., emerging behavior). Fourth, its specification allows for the exploration of a wide range of climate, economic, and policy assumptions. In short, this model serves as a good tool for analyzing how the interplay of complexity (i.e., emergent non-linear behavior) and deep uncertainty can be analyzed through the combination of different computational intelligence

tools. Having said this, the reader should be conscious of the modeling features that fall outside the scope of this study. First, although empirically valuable, the model currently does not consider the possibility of endogenous innovation in the emerging region. Second, international trade and oil prices are also currently outside the scope of this work. Finally, in the optimization setup of this framework, without taxation on fossil fuels, it is not possible to mobilize the resources necessary to fund complementary technology policy for mitigating climate change. In reality, there are many other financial channels through which it would be possible to fund technology-based climate mitigation policies. In the following paragraphs, we describe the most relevant aspects of the model for this analysis¹.

In EDIAM's framework, international climate policy is comprised on nine different elements:

1. Number of years policy intervention is active, starting in 2022: D
2. Carbon tax in the advanced region: τ^A
3. Carbon tax in the emerging region: τ^E
4. Technology subsidy for sustainable energy technologies in the advanced region: h^A
5. Technology subsidy for sustainable energy technologies in the emerging region: h^E
6. R&D subsidy for sustainable energy technologies in the advanced region: q^A
7. R&D subsidy for sustainable energy technologies in the emerging region: q^E
8. Green Climate Fund (GCF) technology subsidy for sustainable energy technologies in the emerging region: h^G
9. GCF R&D subsidy for sustainable energy technologies in the emerging region: q^G .

Formally, the optimal policy intervention is that which maximizes the intertemporal utility of representative consumers in the advanced and emerging regions (Equation 1.1)². which depends both on consumption C (Equation 1.7) and the effects of fossil fuels used in production on temperature rise ΔT (i.e., quality of the environment S , Equations 1.3–1.6), subject to the intertemporal equilibrium conditions of both economies (Equations 1.8–1.13) and to the budget constraint (Equations 1.17, 1.18) in both regions (Acemoglu et al., 2012). The setup of the budget constraints is such that investments on technology-oriented climate action (i.e., technology and R&D subsidies) cannot be greater than the fiscal resources collected through a carbon tax in each region. Cooperation between regions is possible through the use of the GCF,

¹The complete specification of EDIAM, including the derivation of intertemporal dynamics and calibration, can be found in Molina-Perez (2016), specifically Chapter 3, pages 31–59, Appendix A, pages 152–161, and Appendix C, pages 163 and 164. Additionally, we have made publicly available all datasets and programming scripts describing the operationalization of EDIAM's framework in the following github repository: https://github.com/emolinaperez/Ediam_vFrontiers.

²The EDIAM model is specified in continuous form. Yet, it is important to note that the discounted values of utility are computed using discrete time steps after the continuous model is numerically solved.

which redirects resources from the advanced region to the emerging region.

Formally, this is expressed as follows³:

$$Max_{D, \tau^A, \tau^E, h^A, h^E, q^A, q^E, h^G, q^G} \sum_{T_0}^T \frac{1}{(1 + \rho)^t} (u^A + u^E) \quad (1.1)$$

s.t.

$$u^k (C^k, S) = \frac{(\phi(S)C^k)^{1-\sigma}}{1-\sigma} \quad (1.2)$$

$$\phi(S) = \phi(\Delta S) = \frac{(\Delta T_{disaster} - \Delta T(S))^\lambda - \lambda \Delta T_{disaster}^{\lambda-1}}{(1-\lambda)\Delta T_{disaster}^\lambda} \quad (1.3)$$

$$\frac{dS}{dt} = -\xi (Y_f^A + Y_f^E) + \delta \quad (1.4)$$

$$CO_2 = CO_{2|6.0^\circ C} - S \quad (1.5)$$

$$\Delta T = \beta * \ln \left(\frac{CO_2}{CO_{2,0}} \right) \quad (1.6)$$

$$C^k = Y^k - \psi_s \int_0^1 x_{s,i}^k di - \psi_f \int_0^1 x_{f,i}^k di \quad (1.7)$$

$$Y^k = \left(Y_s^k \frac{\varepsilon-1}{\varepsilon} + Y_f^k \frac{\varepsilon-1}{\varepsilon} \right)^{\frac{\varepsilon}{\varepsilon-1}} \quad (1.8)$$

$$Y_j^k = L_j^{k1-\alpha^k} \int_0^1 A_{ji}^{k1-\alpha^k} x_{ji}^{k\alpha^k} di \quad (1.9)$$

$$\frac{Y_s^k}{Y_f^k} = \left(\frac{p_f^k * (1 + \tau^k)}{p_s^k} \right)^\varepsilon \quad (1.10)$$

$$\frac{p_s^k}{p_f^k} = \left(\frac{A_f^k}{A_s^k} \right)^{(1-\alpha^k)} \left(\frac{(1-h^k)\psi_s}{\psi_f} \right)^{\alpha^k} \quad (1.11)$$

$$\frac{\Pi_s^k}{\Pi_f^k} = (1 + q^k) * \frac{\eta_s}{\eta_f} * \frac{1}{(1-h^k)^{\frac{1}{1-\alpha^k}}} * \left(\frac{\psi_f}{\psi_s} \right)^{\frac{\alpha^k}{1-\alpha^k}} * \left(\frac{p_s^k}{p_f^k} \right)^{\frac{1}{1-\alpha^k}} * \frac{L_s^k}{L_f^k} * \frac{A_s^k}{A_f^k} \quad (1.12)$$

$$\theta_j^k = \frac{e^{V(\Pi_j^k)}}{\sum_{j=s}^f e^{V(\Pi_j^k)}} \quad (1.13)$$

$$\psi_j = \psi_{j,0} (x_j^A + x_j^E)^{li} \quad (1.14)$$

$$\frac{dA_j^A}{dt} = \gamma_j (A_j^A) \eta_j \theta_j^A A_j^A \quad (1.15)$$

$$\frac{dA_j^E}{dt} = v_j \gamma_j (A_j^E) \theta_j^E (A_j^A - A_j^E) \quad (1.16)$$

³For clarity, all time subscripts are omitted, but all variables in this optimization set up are dynamic.

$$\sum_{t=D_0}^D \left(h^A \psi_s \int_0^1 x_{s,i}^A di \right) + \left(h^G \psi_s \int_0^1 x_{s,i}^E di \right) + (q^A \eta_s \Pi_s^A) + (q^G \eta_s \Pi_s^E) \leq \sum_{t=D_0}^D \tau^A p_f^A Y_f^A \quad (1.17)$$

$$\sum_{t=D_0}^D \left(h^E \psi_s \int_0^1 x_{s,i}^E di \right) + (q^E \eta_s \Pi_s^E) \leq \sum_{t=D_0}^D \tau^E p_f^E Y_f^E \quad (1.18)$$

where

β : atmosphere's sensitivity to CO₂ emissions (degrees Celsius)

ξ : atmosphere's carbon sink capacity (ppm/BTU/year)

δ : average rate of natural environmental regeneration (dimensionless/year)

ΔT : temperature rise since preindustrial times (degrees Celsius)

CO_{2|6.0 °C}: CO₂ emissions concentration that will result in temperature rise of 6.0°C with respect to preindustrial levels⁴ (ppm)

$\Delta T_{disaster}$: 6.0 (degrees Celsius)

$\phi(S)$: costs of environmental quality degradation (dimensionless)

ε : elasticity of substitution (dimensionless)

ρ : discount rate (dimensionless/year)

p_j^k : primary energy prices of sector "j," region "k" (usd/BTU)

Π_j^k : innovation profitability of sector "j," region "k" (dimensionless)

α^k : proportion of capital income to the total income of the economy in region "k" (dimensionless)

η_j^k : energy technologies propensity to innovation in sector "j," in region "k" (dimensionless/year)

ψ_{ji}^k : unitary cost of production for technology type "i" in sector "j" in region "k" (usd/machine)

L_j^k : share of labor working in sector "j," region "k" (dimensionless)

A_j^k : productivity of sector "j," region "k" (dimensionless)

x_{ji}^k : number of units of technology "i" used in sector "j" in region "k" (machines)

θ_j^k : share of entrepreneurs working in sector "j," region "k" (dimensionless)

γ_j : mean R&D returns to productivity in sector "j" (dimensionless/year)

η_j : innovation propensity of sector "j" (dimensionless/year)

ν_j : probability of successfully imitating/adapting in the emerging region the technologies of sector "j" developed in the advanced region (dimensionless)

T : end of simulation, year 2100

T_0 : initial year of simulation, year 2012

$j \in \{ "s" - \text{sustainable energy} - "f" - \text{fossil energy} - \}$

$k \in \{ "A" - \text{advanced region} - "E" - \text{emerging region} - \}$.

As shown in Equation (1.2), we model consumer preferences through a constant relative risk aversion (CRRA) utility function, which depends both on consumption C and on the quality of the environment S . The parameter σ is the inverse of the intertemporal elasticity of substitution. Equation (1.3) describes the quality of the environment as dependent on temperature rise, which is determined by Equations (1.4)–(1.6), where ΔT represents the increase in average surface global temperature since preindustrial times for a given level of CO₂ atmospheric concentration. The parameter λ controls how quickly the quality of the environment decreases as anthropogenic CO₂ emissions rise. In the same fashion as Acemoglu et al. (2012), the state variable S is a metric of general environmental quality. In this study, this is empirically measured in parts per million (ppm) of atmospheric CO₂ concentrations: the lower the value of S , the higher the environmental quality of the planet. The combination of Equations (1.3)–(1.6) connects this state variable to CO₂ atmospheric concentrations, which in turn allows for internalizing the marginal impact of global fossil energy consumption on consumers' utility.

As shown in Equation (1.7), consumption depends on final production and the cost of technologies used in production. Final production (Equation 1.8) is modeled as a CES aggregate of the two primary energy sources: fossil fuel-based energy (f) and sustainable energy (s). Primary energy production (Equation 1.9) assumes that economic agents use labor and an infinite number of sector-specific technologies "i" for energy production (Acemoglu, 2002), L_j^k represents the labor used in sector "j" $\in \{f, s\}$, A_{ji}^k is the productivity of technology of type "i" used in sector "j", and x_{ji}^k is the number of units of technology type "i" in sector "j" used in production, in region "k." For operationalizing the model, we rely on the same assumption used by Acemoglu et al. (2012): $A_j^k \equiv \int_0^1 A_{ji}^k di$, such that A_j^k is the average productivity of sector "j" in region "k."

The share of production of each energy type "j" (Equation 1.10) depends on the prices of secondary energy types and the carbon tax. Secondary energy prices (Equation 1.11) in turn depend on productivity improvements in both energy sectors, technology costs, and technology subsidies. Technology costs (Equation 1.14) depend on the accumulated number of technologies used in each sector "j" in both regions. The parameter ι_i in this power-law function controls the rate at which experience leads to cost reductions in technology sector "i." The evolution of productivity of section "j" in the advanced region (Equation 1.15) depends on share of entrepreneurs working in this sector, its R&D returns to productivity, and its innovation propensity. For the emerging region (Equation 1.16), we assume that technology entrepreneurs also innovate, but their efforts are targeted toward imitating the existing technologies in the advanced region. The success of these endeavors depends on the ease of transferability of technologies invented in the advanced region. The share of entrepreneurs working on sector "j" (Equation 1.13) determines the sectorial rate of technological progress, which depends on the value " $V(\cdot)$ " that investors assign to the mean profitability of sector "j" in region "k" (Π_j^k). Following the same approach as Train

⁴The model is not defined beyond this limit of temperature rise because such level of temperature rise will result in abrupt and irreversible changes to the global climate system, including events such as the ice sheet collapse, permafrost carbon release, and methane hydrate release (IPCC, 2013).

Kenneth (2003) and Achtenicht et al. (2012), this value function is a deterministic utility component that models economic agents' decisions over competing alternatives in the logistic form expressed in Equation (1.13).

The relative profitability of each sector “j” is described in Equation (1.12). If this ratio is >1, then the majority of research and development is directed toward sustainable energy technologies. In the tradition of Acemoglu (2002) and Acemoglu et al. (2012) framework, Equation (1.12) shows that there are three key forces determining which sector captures the greater share of entrepreneurial activity: (1) the “direct productivity effect” $\frac{A_s^k}{A_f^k}$ incentivizing research in the sector with the more advanced and productive technologies, (2) the “price effect” $\frac{p_s^k}{p_f^k}$ incentivizing research in the energy sector with the higher energy prices, and (3) the market size effect $\frac{L_s^k}{L_f^k}$ pushing R&D toward the sector with the highest market size. In addition to these forces, in the EDIAM modeling framework, two more factors are at play: (1) the “experience effect” $\left(\frac{\psi_f}{\psi_s}\right)^{\frac{\alpha}{1-\alpha}}$ pushing innovative activity toward the sector that more rapidly reduces technological production costs and (2) the “innovation propensity effect” $\frac{\eta_s}{\eta_f}$ incentivizing R&D in the sector that more

rapidly yields new technologies. Note also that the research and technology subsidies also incentivize R&D in sustainable energy technologies. Finally, Equations (1.16) and (1.17) indicate that each region's contribution to the optimal policy should not be greater than the funds collected through the carbon tax.

Table 1 lists the set of policy regimes considered in this study. For each policy, we indicate in which sectors (i.e., carbon tax, technology subsidies, and R&D subsidies) cooperative actions are implemented and in which sectors individual independent actions are carried out. Thus, we model different policy regimes as a mix of individual and cooperative actions across sectors. In total, **Table 1** describes nine different policy regimes. The future without action (FWA) represents the benchmark policy case in which climate policy is not implemented (i.e., laissez-faire economy). The policy regime “P1. I. Carbon Tax [Both]” represents a non-cooperative case in which both regions implement independently climate policy. Policy case “P2. I. Carbon Tax + I.Tech-R&D[Both]” depicts a different non-cooperative policy regime. In this case, the optimal policy response includes independent levels of taxation, technology subsidies, and R&D subsidies for both regions.

Multiple cooperation regimes are described in **Table 1**. For all these policy cases, we assume that regions agree initially on the implementation of a harmonized carbon tax as proposed

TABLE 1 | Description of alternative policy regimes considered.

Policy regime	Independent sectors	Cooperation sectors	Formalism in optimization problem
P0 FWA: Future Without Action	• None	• None	$\tau^A, \tau^E, h^A, h^E, q^A, q^E, h^G, q^G = 0$
P1 I. Carbon Tax [Both]	• Carbon tax	• None	$\tau^A, \tau^E > 0$ $h^A, h^E, q^A, q^E, h^G, q^G = 0$
P2 I. Carbon Tax + I.Tech-R&D[Both]	• Carbon tax • Technology subsidies • R&D subsidies	• None	$\tau^A, \tau^E, h^A, h^E, q^A, q^E > 0$ $h^G, q^G = 0$
P3 H. Carbon Tax + Co-Tech[GCF]+R&D[AR]	• No R&D subsidies in emerging region	• Harmonized carbon tax • Co-funded technology subsidies	$\tau^A = \tau^E > 0$ $h^A, q^A > 0$ $h^E = h^G > 0$ $q^E = q^G = 0$
P4 H. Carbon Tax + Co-Tech[GCF] + I. R&D[Both]	• Independent R&D subsidies	• Harmonized carbon tax • Co-funded technology subsidies	$\tau^A = \tau^E > 0$ $h^A, q^A, q^E > 0$ $h^E = h^G > 0$ $q^G = 0$
P5 H. Carbon Tax + Co-R&D[GCF]+Tech[AR]	• No technology subsidies in emerging region	• Harmonized carbon tax • Co-funded R&D subsidies	$\tau^A = \tau^E > 0$ $h^A, q^A > 0$ $q^E = q^G > 0$ $h^E = h^G = 0$
P6 H. Carbon Tax + Co-R&D[GCF]+I. Tech[Both]	• Independent technology subsidies in emerging region	• Harmonized carbon tax • Co-funded R&D subsidies	$\tau^A = \tau^E > 0$ $h^A, h^E, q^A > 0$ $q^E = q^G > 0$ $h^G = 0$
P7 H. Carbon Tax + Co-Tech-R&D[GCF]	• None	• Harmonized carbon tax • Co-funded R&D subsidies • Co-funded Technology subsidies	$\tau^A = \tau^E > 0$ $h^A, q^A > 0$ $h^E = h^G > 0$ $q^E = q^G > 0$

For each policy regime, it is indicated in which sectors (i.e., carbon tax, technology subsidies, and/or R&D subsidies) cooperative actions are implemented and in which sectors individual independent actions are carried out. Thus, each policy regime can be represented as a mix of individual and cooperative actions across sectors. The set of mathematical restrictions used to represent each policy regime in the optimization framework is noted.

by Nordhaus (2011); therefore, the carbon tax rate is the same across both regions. We also assume that cooperation under the GCF does not have to follow a unique architecture and that it is possible to cooperate in certain sectors, while allowing independent action in others. Policy case “P3: *H. Carbon Tax + Co-Tech[GCF]+R&D[AR]*” considers the case of a harmonized carbon tax across regions and cooperation in co-funded technology subsidies under GCF. However, in this case, independent R&D subsidies are only implemented in the advanced region. Policy “P4: *H. Carbon Tax + Co-Tech[GCF] + I. R&D[Both]*” expands on the latter case by considering that independent R&D subsidies are implemented in both regions.

Policy “P5: *H. Carbon Tax + Co-R&D[GCF]+Tech[AR]*” includes the implementation of a harmonized carbon tax in both regions, co-funded R&D subsidies under the GCF and independent technology subsidies in the advanced region. Policy regime “P6: *H. Carbon Tax + Co-R&D[GCF]+I.Tech[Both]*” expands policy case P5 by allowing for the implementation of independent technology subsidies in both regions. Finally, policy regime “P7: *H. Carbon Tax + Co-Tech-R&D[GCF]*” considers the case in which in addition to a harmonized carbon tax, cooperation under the GCF includes co-funded R&D subsidies and technology subsidies.

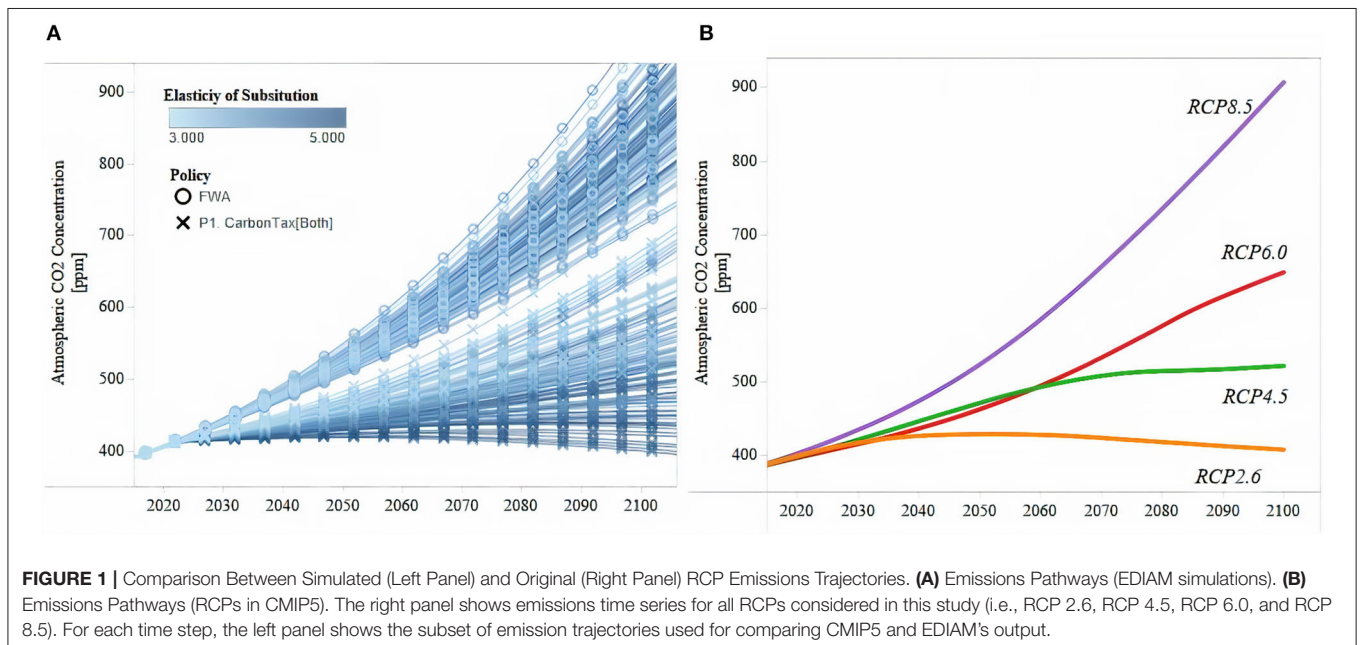
Uncertain Stressors Across Spheres

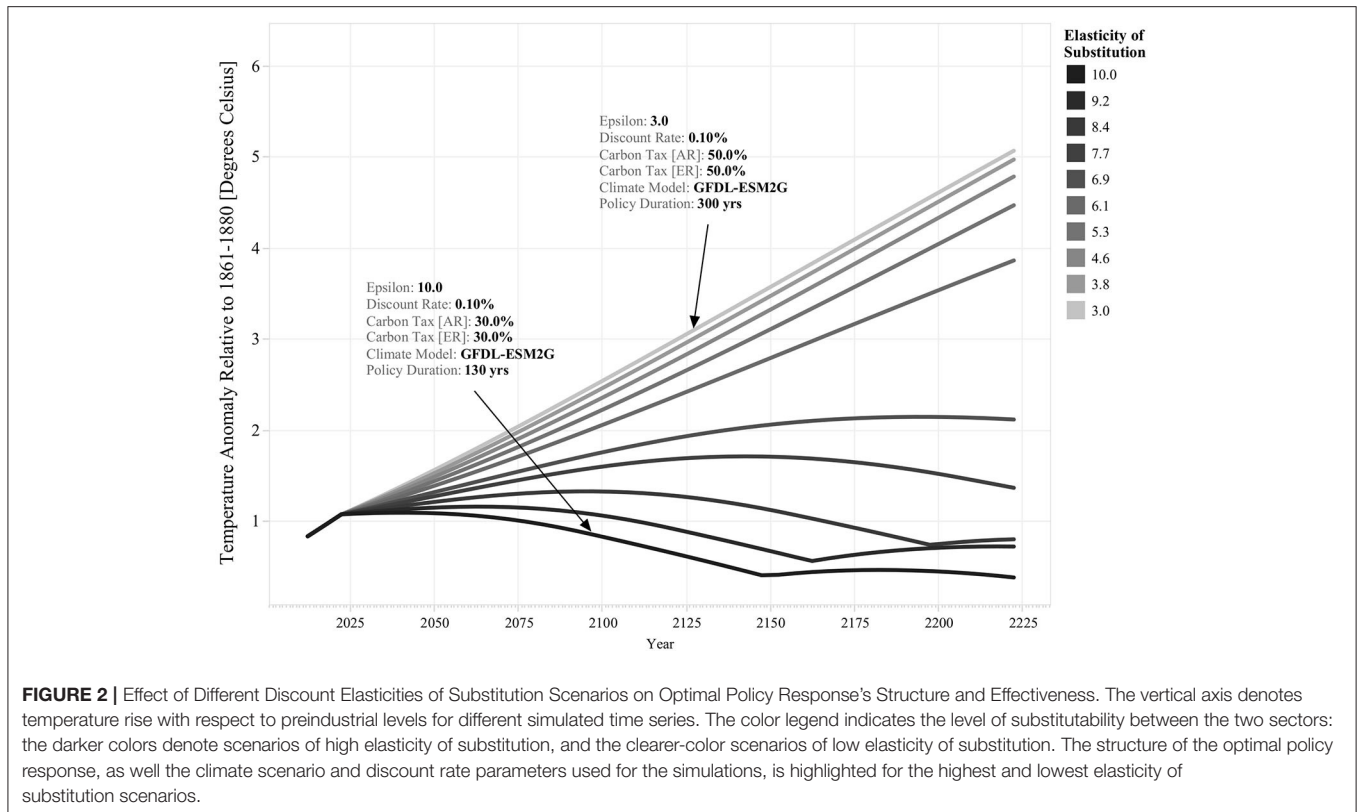
To analyze the performance of different policy regimes across uncertainty, we focus on four uncertain stressors, affecting two spheres: (1) the elasticity of substitution between fossil and sustainable energy inputs in production, and the economic agents’ discount rate, impacting the economic sphere and (2) climate sensitivity to GHG emissions and the capacity of atmospheric carbon sinks affecting the ecological sphere.

Thus, there are two types of elements in our analysis: (1) policy regimes that describe different sectorial interventions and cooperation schemes between regions and (2) scenarios which describe unique parameter combinations of economic and climatic variables.

This framework allows us to explore uncertainty in more detail by generating an ample set of emission trajectories through variations of economic parameters and policy regimes. For example, **Figure 1** compares a subset of simulated emission pathways that vary the elasticity of the substitution parameter (ϵ , Equation 1.8) for two policy regimes, against the four Representative Concentration Pathways (RCPs) included in the CMIP5 dataset. It is possible to see that the range of variation produced with these simulations is similar to that captured by the four RCPs included in CMIP5. This feature is important for this analysis because as discussed in section Machine Learning Algorithms for Identifying Decision-Relevant Conditions, by considering such a disaggregated set of variation, it is possible to identify with higher precision vulnerability thresholds.

The elasticity of substitution is an important parameter in the economic sphere because it describes the extent to which sustainable energy technologies can be used to substitute the functions of fossil energy technologies in secondary energy production. The results of Acemoglu et al. (2012) have spurred interest among empirical researchers on estimating more accurately the potential level of substitution between the two sectors. At present, initial empirical results show that the short- and long-term values of the elasticity of substitution are likely to be closer to the low substitution case considered in Acemoglu et al. (2012), but more importantly, these initial results show that the strength of the substitution effect in the long term is highly uncertain. For instance, Papageorgiou et al. (2013) use



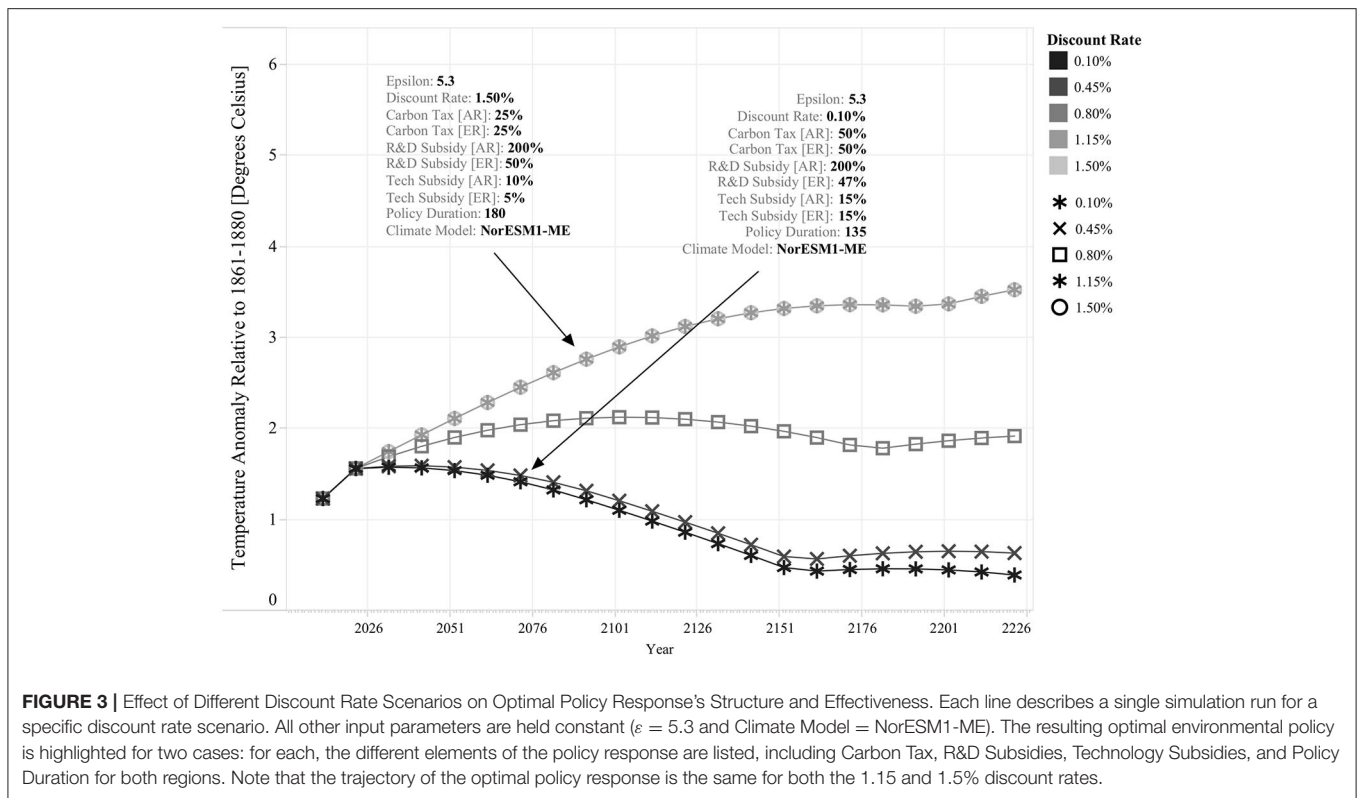


cross-country sectoral energy data and nested CES production functions to estimate this parameter. They find evidence that the elasticity of substitution in the short term is more likely to be in the low substitution range ($\epsilon = 3$) of Acemoglu et al. (2012) study, but in the long term it is plausible that this parameter falls in the high-range values. Another study by Pottier et al. (2014) argues that the elasticity of substitution between sustainable energy and fossil energy is also likely to be in the low substitution range ($\epsilon = 3$), perhaps even below one ($\epsilon < 1$). They argue that this is the case mainly because capital stocks for most of the energy system last for many decades, and this delays substitution away from fossil energy. However, the authors consider that in the long term, as innovation broadens the range of technological possibilities, it is plausible that all energy sources will be fairly substitutable. These results and the debate among researchers on this topic support the notion that the elasticity of substitution between sustainable and fossil energy is a deeply uncertain parameter. These empirical findings show that the current state of science does not provide sufficient and adequate evidence to estimate accurately this parameter. They also suggest that in the long term a wide range of values is plausible.

We explore the implications of varying levels of substitutability between the two sectors by considering 10 different scenarios for this parameter. **Figure 2** lists the different elasticity of substitution scenarios considered in this analysis and exemplifies their effect on temperature rise stabilization. It shows that the 10 scenarios considered for the elasticity of substitution can result in substantially different outcomes. For

instance, it shows that for three high levels of substitutability scenarios, $\epsilon = 10.0$, $\epsilon = 9.2$, and $\epsilon = 8.4$, it is possible to induce a full self-reinforcing transition away from fossil energy before the end of the simulation runs (i.e., policy duration < 300 years). In contrast, for the low elasticity of substitution scenarios, $\epsilon = 3.0$, $\epsilon = 3.8$, $\epsilon = 4.6$, and $\epsilon = 5.3$, it is necessary to sustain policy intervention (i.e., harmonized carbon tax in both regions) during the entire simulation at a high level (i.e., 50%) to delay temperature rise. This shows that the cost and effectiveness of policy intervention is closely linked to the degree of substitutability between the fossil and sustainable energy sectors. The less substitutable these sectors are, the more effort is required to induce a successful transition toward sustainable energy and the decarbonization of secondary energy production in both regions.

The discount rate is a mathematical formalism that helps us express future costs and gains at today's equivalent value. In the context of climate change, this parameter attempts to describe how societies of today value the environmental and economic outcomes of the future. Controversy over the proper value of the discount rate lies at the heart of many of the debates associated with climate change policy. It should not be a surprise that studies that use different discounting values reach different conclusions regarding the structure of the optimal environmental policy required to stabilize global temperature rise. This debate is best exemplified by Nordhaus and Stern's research on the level of carbon taxation needed to keep temperature rise at sustainable levels (Acemoglu et al., 2012). In short, Nordhaus,



using a discount rate of 1.50% per year, finds that an initial small carbon tax that increases over time would guarantee that temperature rise will be kept below three degrees Celsius in the long term, while Stern, using a discount rate of 0.10% per year, argues that a higher initial carbon tax is needed to achieve temperature rise stabilization sooner and avoid future significant damage from climate change. This disagreement among climate experts is evidence of the deep uncertainty associated with the discount rate.

In this analysis, we explore this uncertain stressor by considering a diverse set of discount rate scenarios. To develop these scenarios, we assume that the maximum value that this parameter can take is the one proposed by Nordhaus (i.e., 1.5% per year) and that the minimum value is the one proposed by Stern (i.e., 0.10% per year). However, we also consider three more possibilities in between to explore in more detail the role of varying levels of discounting on the structure of the optimal policy response.

Figure 3 lists the five discount rate scenarios considered in this analysis. This exercise provides an illustrative example of the discount rate's role in determining the structure of the policy response. By comparing the optimal policy response across the Stern (i.e., 0.10% per year) and Nordhaus (i.e., 1.5% per year) limits, it is possible to see that in the first case policy intervention is more decisive across both regions than policy intervention in the second case. For instance, the policy response with the 0.10% per year discount rate uses higher levels of carbon taxation and technology subsidies in both regions. As a result, the

environmental outcomes are also significantly different; for the 0.10% discount rate, temperature rise is kept below two degrees Celsius throughout the entire simulation, while for the 1.15% discount rate, temperature rise continues for over a century until it is stabilized at $\sim 3^\circ\text{C}$. In this case, the cost of policy intervention is higher for the 0.10% discount rate, but it is important to note that in comparison to the 1.5% discount rate policy, this policy requires to be implemented during a shorter period of time (i.e., 135 vs. 180 years); thus, under alternative climate conditions, it is also feasible that both policies display similar intervention costs, or that in fact, the 0.10% discount rate policy becomes cheaper.

The uncertainty associated with the speed of temperature rise is associated with the limitations of our understanding of the global climate system. Each general circulation model used by the IPCC and included in the CMIP5 ensemble uses different assumptions and parameter values to describe the atmospheric changes resulting in growing anthropogenic GHG emissions, and, as a result, the magnitude of the estimated changes varies greatly among different modeling groups. In this respect, one of the features of EDIAM is that it uses 12 GCMs included in the CMIP5 data ensemble to calibrate the parameters ξ , δ , β , and S_0 in Equations (1.4) and (1.6). Thus, in EDIAM, GCMs are described as unique combinations of climate sensitivity of GHG (β) and the capacity of the atmospheric carbon sink (δ , S_0) as listed in **Table 2**.

Figure 4 provides an illustrative example of how different GCMs may lead to a different structure of the optimal environmental policy. It shows that for a GCM that displays

higher climate sensitivity, such as *MIROC-ESM-CHEM*, it is possible that under certain circumstances, the optimal policy uses a higher mix of carbon taxes, research subsidies, and technology subsidies than in the case of a GCM that displays lower climate sensitivity, like *NorESM1-M*. It also shows that the environmental outcomes between both scenarios are different: in this case, for

both simulation runs temperature rise is successfully mitigated, but this occurs at a higher level for climate scenario *MIROC-ESM-CHEM* than for scenario *NorESM1-M*. It also shows that the cost of policy intervention is unambiguously higher for climate scenario *MIROC-ESM-CHEM* because although the rate of carbon taxation is smaller in climate scenario *NorESM1-M*, policy intervention lasts longer in the latter case. Evidently, these results can change when combined with other uncertainties, yet it offers an illustrative example of the interplay between the optimal policy response and the different climate scenarios.

TABLE 2 | Estimated climate parameters using CMIP5 GCM models.

Climate scenario	β	ξ	δ	S_0
MIROC-ESM-CHEM	6.13	0.010	0.00278	590
GFDL-CM3	6.11	0.010	0.00259	635
MIROC-ESM	5.93	0.010	0.00260	633
bcc-csm1-1	5.00	0.010	0.00182	916
MPI-ESM-LR	4.67	0.010	0.00161	1,042
MPI-ESM-MR	4.67	0.010	0.00161	1,045
NorESM1-ME	4.34	0.010	0.00136	1,236
MRI-ESM1	4.26	0.010	0.00130	1,294
NorESM1-M	4.13	0.010	0.00119	1,415
MIROC5	4.12	0.010	0.00119	1,417
GFDL-ESM2M	3.29	0.010	0.00071	2,403
GFDL-ESM2G	3.19	0.010	0.00063	2,695

The table lists the estimated parameters for the 12 CMIP5 climate models included in this study; the parameters of Equations (1.4) and (1.6) are listed for each climate model. These parameters are estimated using CO₂ emission levels' variation across representative concentration pathways (RCPs) for each of the climate models in an autoregressive model.

USING MODELS DIFFERENTLY THROUGH COMPUTATIONAL EXPERIMENTATION

Considering Multiple Dimensions of Merit

Sustainability challenges often deal with multiple spheres (e.g., economic, ecological, technological) (Liu et al., 2013; Hull et al., 2015); as a result, sustainability studies need to deal with multiple, and often, opposing measures of merit. Climate change mitigation offers a clear example of this as it requires the consideration of different metrics to evaluate and compare the performance of competing policy proposals. In this study, we focus primarily on the outcome that policy intervention has on economic and environmental conditions by the end of the century. This aligns the scope of this work to discussions associated with the end-of-the-century temperature rise and emission stabilization targets.

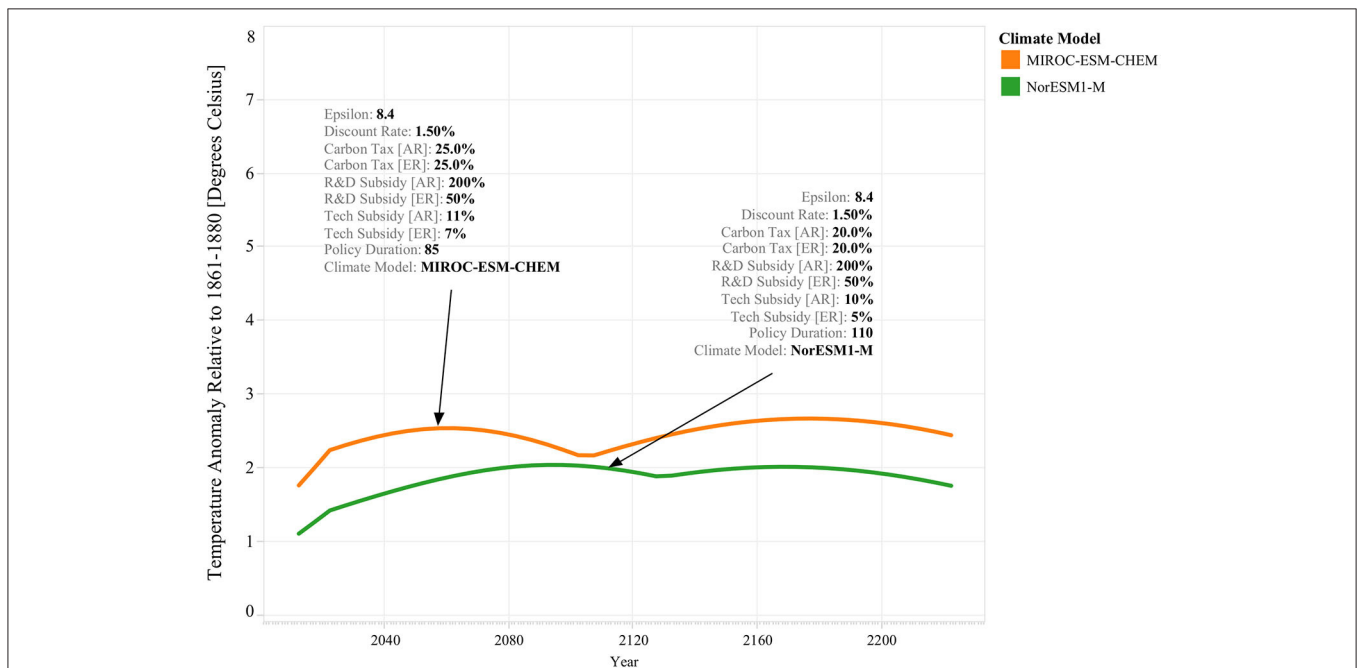


FIGURE 4 | Effect of Different Climate Scenarios on Optimal Policy Response's Structure and Effectiveness. This figure shows temperature rise time series for two simulation experiments. Both simulations used the same parameter values for the elasticity of substitution (i.e., 8.4) and the discount rate (i.e., 1.50% per year) but are run for different climate scenarios: *MIROC-ESM-CHEM* (i.e., orange line) and *NorESM1-M* (i.e., green). For each simulation, the pointing arrows indicate the resulting optimal policy as a combination of carbon taxes, research subsidies, and technology subsidies across both regions.

From an economic perspective, we estimate the cost of policy intervention by comparing consumption levels across the policy intervention case and the laissez-faire economy. Then, the higher the reduction in consumption compared to the laissez-faire economy, the higher the costs of policy intervention. From an environmental perspective, we consider two metrics: the end-of-the century temperature rise level and end-of-century CO₂ atmospheric concentrations. The first metric is useful for comparing policies in terms of the temperature levels that are plausible with its implementation. The second metric is useful to analyze whether or not a policy stabilizes CO₂ emissions such that temperature permanently stops rising. We make this distinction because maintaining temperature rise below a certain threshold (e.g., two degrees Celsius) does not entail that atmospheric CO₂ concentrations are also stabilized. Without stabilization, if climate policy is lifted, temperatures will continue rising.

Experimental Design and Case Generation

We use the elements outlined in the previous sections to conduct several simulation experiments. The experimental design includes a full-factorial sampling design across different EDIAM's parameters; this includes

- 12 climate scenarios
- 10 elasticity of substitution scenarios
- 5 discount rate scenarios.

We considered all possible combinations of these uncertain exogenous factors for developing individual model parametrizations, which yields a total of 600 cases. **Table 3** summarizes the scope of the experimental design of this study using the XLRM framework developed by Lempert (2003), while emphasizing that we are dealing specifically with uncertain stressors in the context of sustainability (i.e., XSLRM).

MACHINE LEARNING ALGORITHMS FOR IDENTIFYING DECISION-RELEVANT CONDITIONS

Experimental Datasets and Results

Figure 5 describes the sequence of steps we implemented to produce the datasets used for the analysis described in this section. For each of the steps in the process, this figure indicates the method and general characteristics of the datasets produced. The experimental design consisted of 5,400 optimization runs across 600 parametrizations that vary climate parameters, elasticity of substitution, and the discount rate. The optimization runs estimate the optimal policy response for each of the parametrizations, considering the restrictions of the different policy regimes, using Byrd et al. (1995) "L-BFGS-B" method for constraint optimization. On average, it takes 10,000 simulation runs to converge on a solution for the optimization problem. Thus, in total, the results described in the following sections required ~54 million simulation runs.

Four datasets are relevant for this sequence of steps. The experimental design dataset describes how the combination

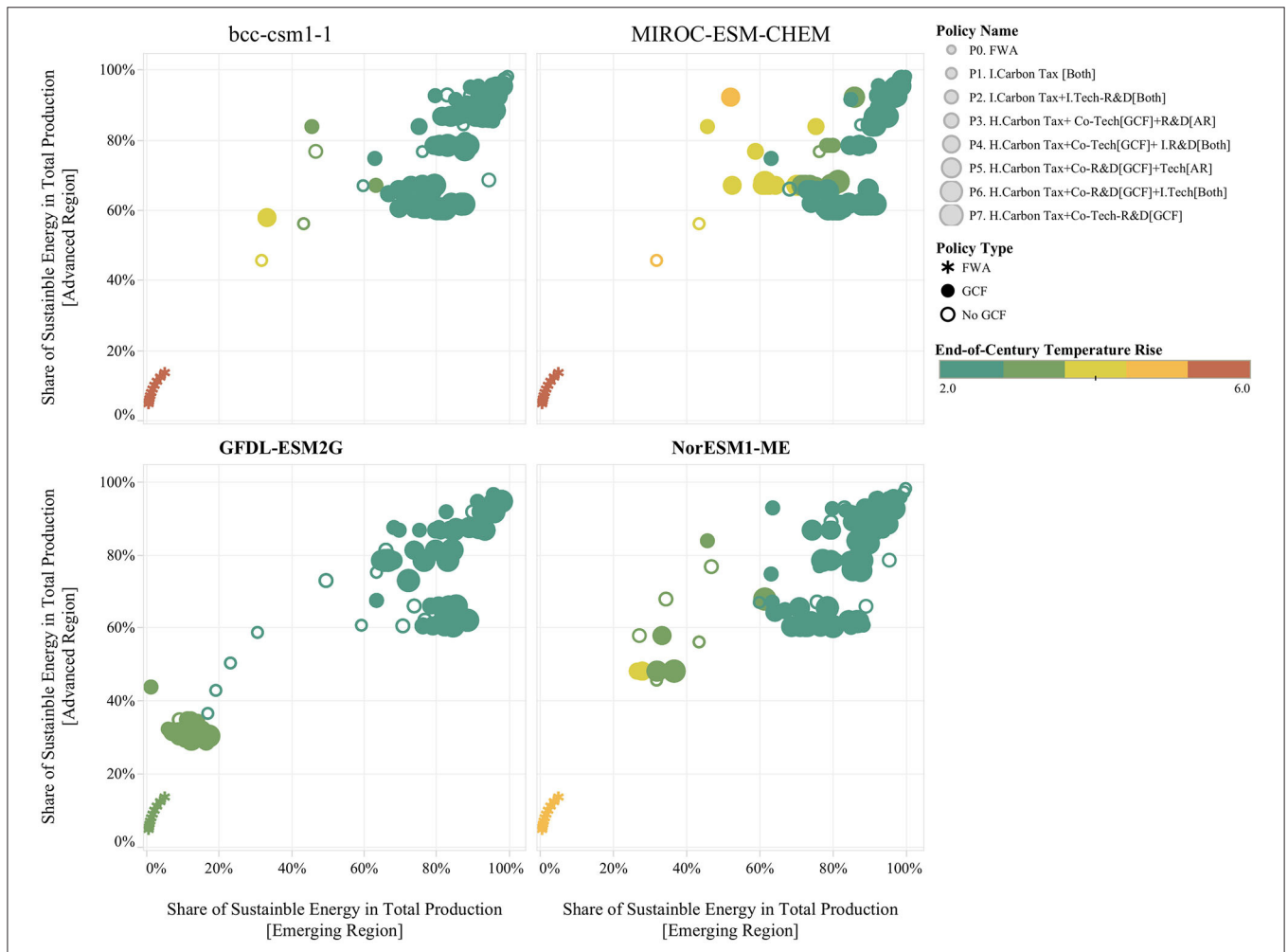
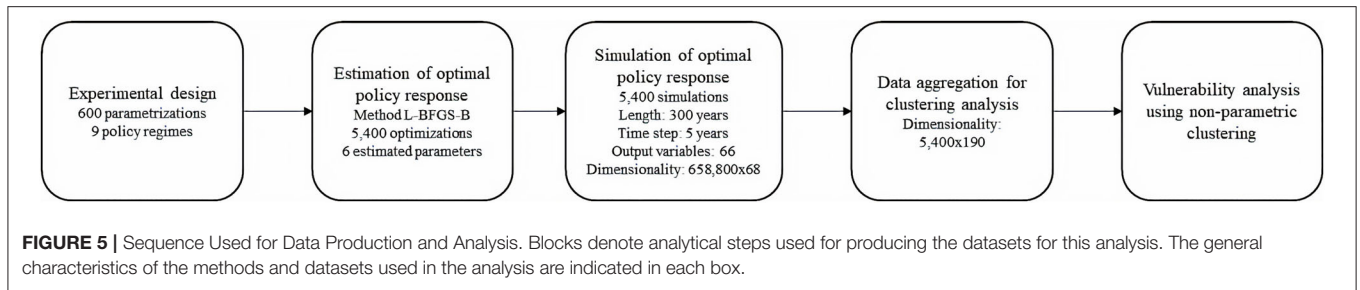
TABLE 3 | XSLRM summary of experimental design.

Uncertain stressors (XS)	Policy levers (L)
Climate uncertainty: <ul style="list-style-type: none"> • 12 Climate scenarios 	<ul style="list-style-type: none"> • P0. FWA (Future Without Action) • P1. I. Carbon Tax [Both]
Economic uncertainty: <ul style="list-style-type: none"> • 10 elasticity of substitution scenarios • 5 discount rate scenarios 	<ul style="list-style-type: none"> • P2. I. Carbon Tax + I.Tech-R&D[Both] • P3. H. Carbon Tax + Co-Tech[GCF]+R&D[AR] • P4. H. Carbon Tax + Co-Tech[GCF] + I. R&D[Both] • P5. H. Carbon Tax + Co-R&D[GCF]+Tech[AR] • P6. H. Carbon Tax + Co-R&D[GCF]+I. Tech[Both] • P7. H. Carbon Tax + Co-Tech-R&D[GCF]
System relationships (R)	Metrics (M)
<ul style="list-style-type: none"> • Exploratory dynamic integrated assessment model (EDIAM) 	<ul style="list-style-type: none"> • End-of-century temperature rise • Stabilization of GHG emissions • Economic costs of policy intervention

The main components of the exploratory analysis are grouped according to four different categories: (1) the deep uncertainty scenario taken into account (i.e., 12 climate scenarios, 10 Elasticity of Substitution Scenarios, and 5 Discount Rate Scenarios), (2) the policy regimes analyzed (i.e., 8 different policy regimes), (3) the system relationship that links actions to consequences (i.e., EDIAM model), and (4) the metrics considered to analyze the performance of different policies.

of climate and economic parameters vary across the different optimization runs. In terms of its cardinality, there are 600 unique combinations of parameters in this dataset, identified by unique future ids, which are combined with the 9 policy regimes, indicated by a unique policy id. The optimal policy response dataset describes for each of the 5,400 runs the combination of policy parameters that solves the optimization problem described in section Virtual Laboratories and Policy Regimes; each of these optimal vectors is unique, since each of the estimated variables is continuous (e.g., carbon tax rates, subsidy rates, and R&D intensities). The simulation dataset describes the dynamic behavior of the system under these 5,400 optimal policy vectors using 66 output variables of the EDIAM model. Finally, in the scenario discovery dataset, we aggregate simulation results by summarizing the dynamic behavior of each run using an expanded set of variables that compare absolute and relative behavior across regions and sectors. For instance, by comparing end-time technological progress with respect to initial conditions, technological progress in competing sectors within regions, and technological progress across regions.

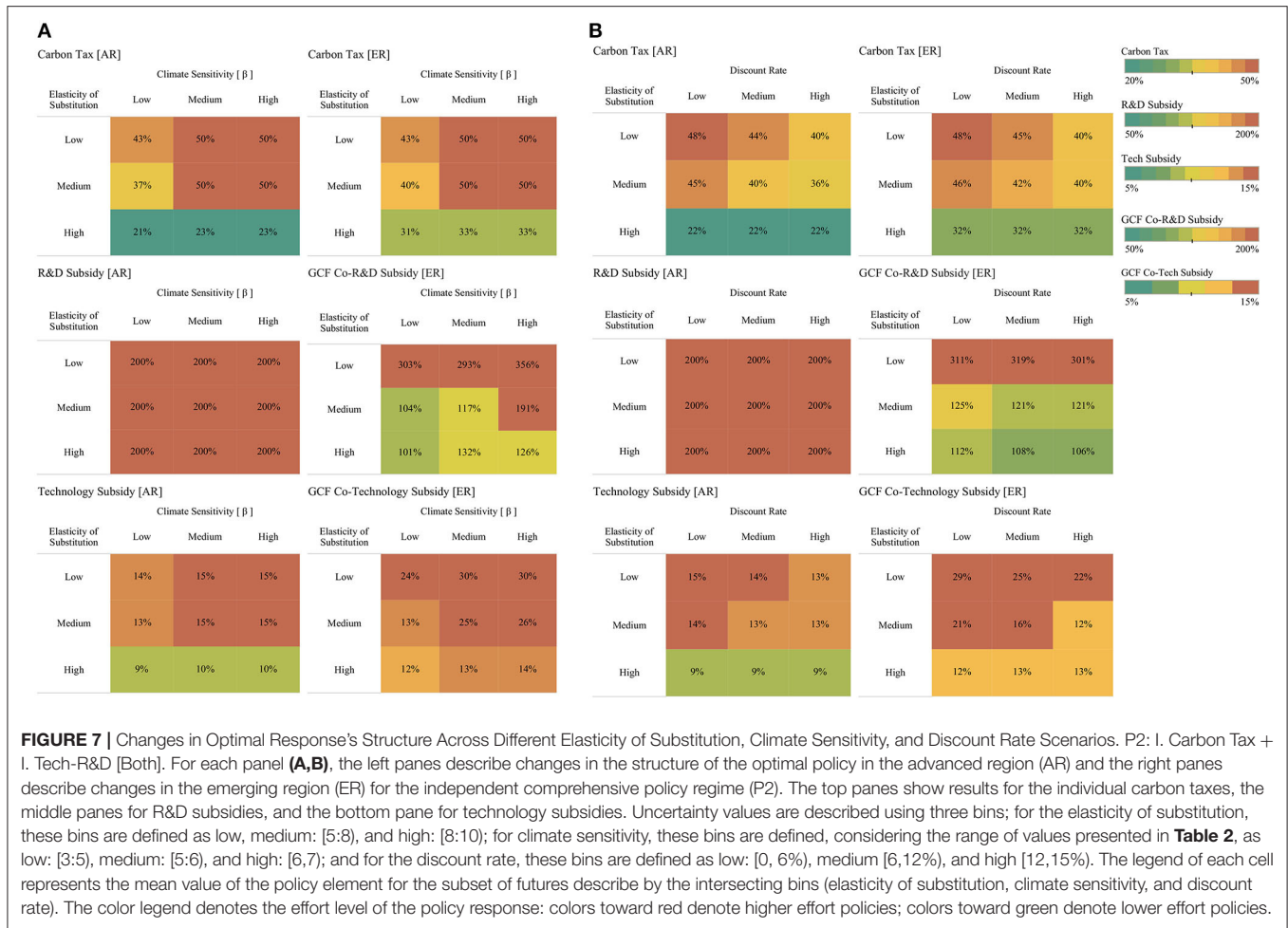
Results from the simulation runs generated by the experimental design are shown in **Figure 6**. These results are useful for highlighting some of the features of the system and of the policy response. The figure shows that there is ample variation with respect to the penetration levels of sustainable energy that can be achieved through the various policy regimes. It is possible to see that the FWA (i.e., laissez-faire economy) results in limited penetration of sustainable energy across both regions, and as a result, end-of-century temperature rise levels are close to the environmental limit (i.e., 6°C). **Figure 6** also reveals that the best environmental outcomes are concentrated in the upper right corner of these panes. These futures represent scenarios in which the policy response induces a successful transition toward



sustainable energy across both regions. It is possible to see that the non-GCF policies (i.e., P1 and P2) can achieve similar levels of penetration of sustainable energy in both regions than GCF-based policies (i.e., P3–P7). These results also show that policy performance varies across GCMs in terms

of the penetration of sustainable energy and the resulting temperature rise.

The structure of the optimal environmental policy varies across the uncertainty space in order to meet the climate policy targets described. This variation in the structure of the



optimal response has important implications for policy design and exemplifies the richness of the experimental results. For example, **Figure 7A** shows for policy P2: I. Carbon Tax + I. Tech-R&D [Both] how the optimal response changes across two parameters: the elasticity of substitution and climate sensitivity to GHG. The left panes describe changes in the structure of the optimal policy in the advanced region (AR), and the right panes describe changes in the emerging region (ER). The top panes show results for the individual carbon taxes, the middle panes for R&D subsidies, and the bottom pane for technology subsidies. The results presented in this figure show that the structure of the optimal policy is very sensitive to the combined effect of the elasticity of substitution and climate sensitivity: the higher the climate sensitivity and the lower the elasticity of substitution, then the higher the effort of the optimal policy response.

The discount rate is another important factor that influences the structure of the optimal policy response. **Figure 7B** describes changes in the structure of the optimal policy across different scenarios of the elasticity of substitution and the discount rate. As expected, it shows that the strength of the policy response increases as the discount rate diminishes. However, in this case it is possible to see that as the elasticity of substitution increases, the influence of the discount rate in the structure of the optimal

policy diminishes. For high elasticity of substitution scenarios, it is possible to see that the structure of the optimal policy is insensitive to changes in the discount rate. These results highlight the importance of regional differences in defining the structure of optimal environmental regulation. It is possible to see that in the emerging region carbon taxation is always equal or higher than carbon taxation in the advanced region. In contrast, the technology policy elements of optimal environmental regulation are higher in the advanced region than in the emerging region. Since technologies are developed in the advanced region, then the optimal policy prioritizes accelerating technology development over taxation in this region, while in the emerging region, higher taxation creates a strong market niche for sustainable energy, which is used more effectively by R&D and technology subsidies that accelerate the technological catching-up process.

A similar analysis for policy regime P7 “H. Carbon Tax + Co-Tech-R&D[GCF]” shows that under the GCF the level of carbon taxation reduces for both regions compared to the level of taxation in the non-cooperative policy regime (i.e., P2). Additionally, the optimal level of effort in R&D and technology subsidies in the emerging region is on average higher than the optimal level of effort in the non-cooperative policy regime. This indicates that under the GCF, it is feasible for the emerging region

to make higher investments in R&D and technology subsidies and reduce the rate of taxation. Similarly, for the advanced region, these results show that it is possible to reduce the level of carbon taxation by co-funding R&D and technology subsidies in the emerging region. Finally, the results show that in the most adverse scenarios under the GCF (i.e., low elasticity of substitution and high climate sensitivity), optimal environmental regulation requires higher R&D and technology subsidies in the emerging region than in the advanced region.

Machine Learning Algorithms for Describing Vulnerability Conditions

The previous section describes general characteristics of the experimental datasets and insights of the computational experiment. Yet, these results do not provide a systemic understanding of how the interaction of the set of stressors considered in the experiment affect the structure and effectiveness of optimal climate policy response under uncertainty. To address this, we follow two steps. First, we classify experimental outcomes with respect to whether or not specific policy objectives are met. Second, we use non-parametric clustering analysis for understanding the combination of factors that lead to meeting these objectives. We consider an outcome is not vulnerable when the temperature target (i.e., 2°C) and/or the stabilization targets are met. This suggests that there are two outcome types of interest in this experiment:

1. Simulations in which the 2°C end-of-century temperature rise target is met
2. Simulations in which the 2°C end-of-century temperature rise target and CO₂ stabilization are met.

Table 4 presents the performance statistics of different policy regimes across the 600 parametrization cases considered for these two outcome types. As expected, the FWA does not meet any of the climate change objectives. It also shows that for the independent carbon tax policy (i.e., P1) in the majority of simulations, it is possible to keep the temperature rise below 2°C, but in none of these cases is this policy able to stabilize CO₂ emissions. This shows that this policy is effective in delaying temperature rise but is less effective at inducing successful decarbonization across regions. In contrast, policies that complement carbon taxes with R&D and technology subsidies are able to meet the CO₂ stabilization targets in a higher number of futures. It is possible to see that the stabilization targets are met in less than one third of the futures considered. In this respect, some of the GCF-based policies (i.e., P4 and P7) are slightly more effective than the non-GCF policy (i.e., P2) in meeting the stabilization target.

For the second step, we use the algorithm PRIM (Patient Rule Induction Method) (Friedman and Fisher, 1999), a non-parametric bump hunting classification algorithm, to quantitatively describe vulnerability condition of different policies. In particular, we use PRIM in the context of the scenario discovery method developed by Bryant and Lempert (2010). Thus, for each policy regime, we classify simulation outcomes into two cases of interest (I_s): (1) cases in which the policy

TABLE 4 | Performance of optimal policy response across different policy regimes.

Policy name	Number (percentage) of futures meeting the end-of-century climate policy target	
	Temperature rise below 2°C	
	CO ₂ stabilization achieved	CO ₂ stabilization not achieved
P0. FWA	0 (0)	0 (0.0)
P1. I. Carbon Tax [Both]	0 (0)	375 (62.5)
P2. I. Carbon Tax + I. Tech-R&D[Both]	153 (25.5)	398 (66.3)
P3. H. Carbon Tax + Co-Tech[GCF] + R&D[AR]	130 (21.7)	344 (57.3)
P4. H. Carbon Tax + Co-Tech[GCF] + I. R&D[Both]	153 (25.5)	391 (65.2)
P5. H. Carbon Tax + Co-R&D[GCF] + Tech[AR]	130 (21.7)	395 (65.8)
P6. H. Carbon Tax + Co-R&D[GCF] + I. Tech[Both]	145 (24.2)	415 (69.2)
P7. H. Carbon Tax + Co-Tech-R&D[GCF]	165 (27.5)	402 (67.0)

The table summarizes the performance of each policy across the 600 parametrizations considered for the four outcome types. The numbers (percentage) of parametrization meeting the different end-of-century climate policy targets are listed under each column.

target is met and (2) cases in which the policy target is not met. Then, PRIM is used to parse the simulation database into concise clusters that describe dimensional conditions under which policies do not meet targets. This is done through the estimation of recursive peeling trajectories, as class types often require more than one cluster to be fully described. This implies that once an initial cluster is chosen, the algorithm removes all the data points from the dataset inside the first cluster and replicates the peeling/pasting process with the remaining data.

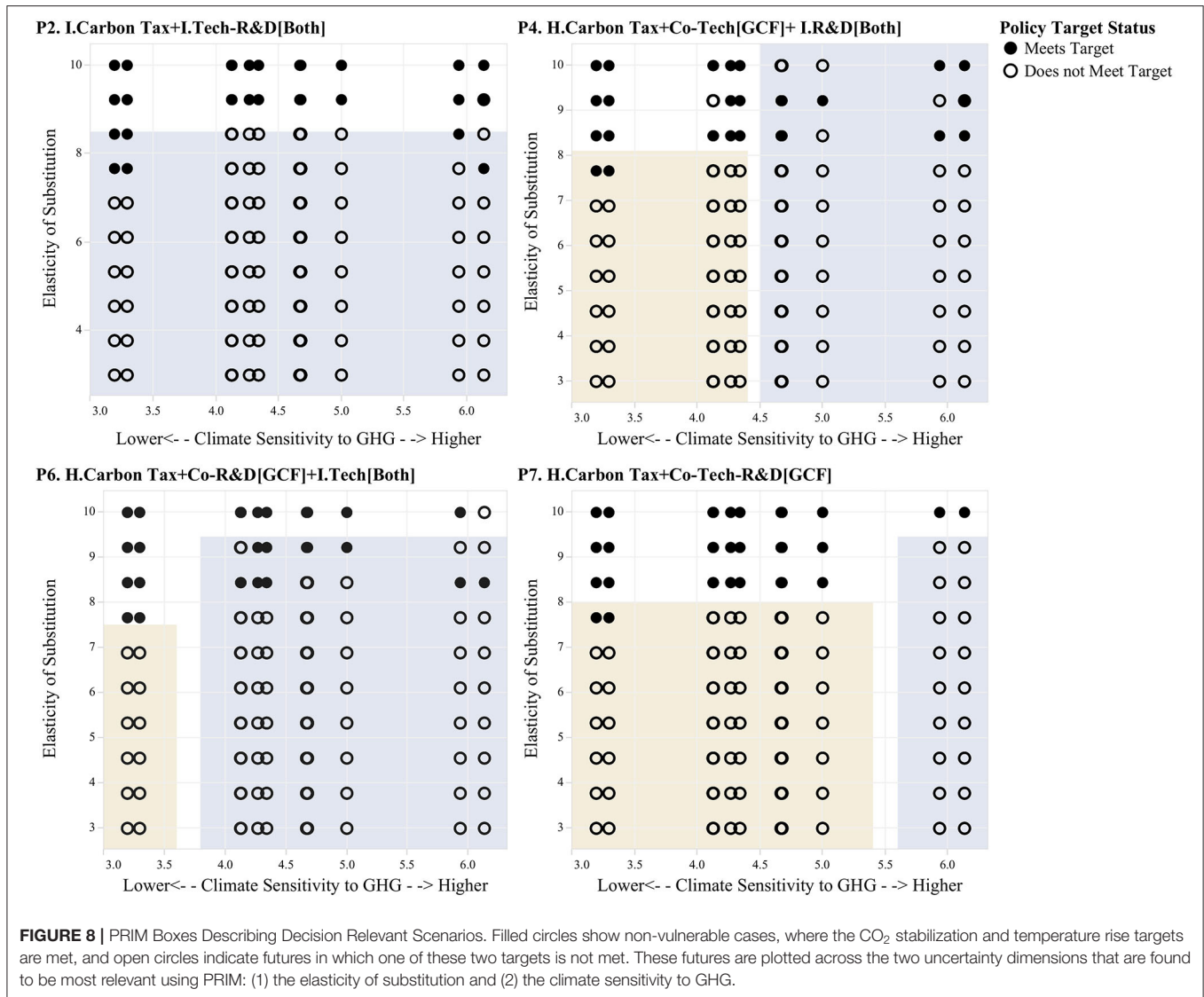
Two statistical measures are used to describe the suitability of a decision relevant cluster. Coverage (Equation 2) measures how completely the cases defined by cluster B cover the cases of interest (I_s); in this study, this is the percent of total vulnerable cases that are captured by the cluster. Density (Equation 3) measures the purity of the scenarios; in this study, this is the percent of cases within the cluster that are vulnerable. Interpretability of these cluster is an important subjective measure; generally, the fewer dimensions used by the cluster, the higher its suitability for the analysis.

$$Coverage = \frac{\sum_{x_i \in B} y_i'}{\sum_{x_i \in I_s} y_i'} \quad (2)$$

$$Density = \frac{\sum_{x_i \in B} y_i'}{\sum_{x_i \in B} 1} \quad (3)$$

where $y_i' = 1$ if $x_i \in I_s$ and $y_i' = 0$ otherwise.

We first use scenario discovery to understand the cases in which end-of-century CO₂ stabilization at 2°C targets is not met. These are futures in which CO₂ stabilization is not achieved



and in which end-of-century temperature rise is above 2°C. **Figure 8** shows the results of this clustering analysis. The figure shows a series of scatter plots of all futures for different policy regimes. Filled circles show non-vulnerable cases, where the CO₂ stabilization and temperature rise targets are met, and open circles indicate futures in which one of these two targets is not met. These futures are plotted across the two uncertainty dimensions that are found to be most relevant using PRIM: (1) the elasticity of substitution and (2) the climate sensitivity to GHG. High values of the elasticity of substitution describe scenarios in which the technologies across sectors are highly substitutable, which are more favorable for climate policy. Low values of the elasticity of substitution denote scenarios in which sectors are less substitutable, which makes it harder to move away from fossil energy. For the case of climate sensitivity, high values describe climate scenarios in which global temperature rises rapidly with growing CO₂, thus making it harder to keep temperature levels below the 2°C target. Low values are

associated with climate scenarios for which global temperature rises less abruptly with growing CO₂ emissions. Finally, the shaded regions highlighted in yellow and blue were selected using scenario discovery to describe these sets of vulnerable futures. **Table 5** provides a detailed description of the boundary conditions of each scenario box, as well as the corresponding coverage and density statistics that describe to which extend these scenario boxes adequately capture the vulnerable conditions of each policy.

The results presented in **Figure 8** and **Table 5** show that the vulnerability region varies slightly across the different environmental policy regimes. For the independent comprehensive policy (“P2 I. Carbon Tax+I. Tech-R&D[Both]”), the vulnerability region is defined solely by the elasticity of substitution. The optimal policy under this regime fails to meet the stabilization target in all scenarios that do not display a high elasticity of substitution. For the other three policy regimes, the vulnerability region is described by both the elasticity of

TABLE 5 | Scenario discovery analysis summary results for stabilization target.

Policy name	Scenario box	Scenario description	Coverage	Density
P2. I. Carbon Tax + I. Tech-R&D[Both]	Box1	<ul style="list-style-type: none"> Elasticity of substitution < 9.0 	99% (445/447)	93% (413/447)
P4. H. Carbon Tax + Co-Tech[GCF] + I. R&D[Both]	Box1	<ul style="list-style-type: none"> Climate sensitivity to GHG > 4.5 	53% (237/447)	80% (190/237)
	Box2	<ul style="list-style-type: none"> Elasticity of substitution < 8.0 Climate sensitivity to GHG < 4.5 	45% (202/447)	95% (182/202)
P6. H. Carbon Tax + Co-R&D[GCF] + I. Tech[Both]	Box1	<ul style="list-style-type: none"> Elasticity of substitution < 9.5 Climate Sensitivity to GHG > 4.0 	86% (392/455)	87% (341/392)
	Box2	<ul style="list-style-type: none"> Elasticity of substitution < 7.6 Climate sensitivity to GHG < 4.0 	13% (60/455)	100% (60/60)
P7. H. Carbon Tax + Co-Tech-R&D[GCF]	Box1	<ul style="list-style-type: none"> Elasticity of substitution < 9.5 Climate sensitivity to GHG > 5.5 	30% (130/435)	97% (126/435)
	Box2	<ul style="list-style-type: none"> Elasticity of substitution < 8.0 Climate sensitivity to GHG > 5.5 	70% (305/435)	97% (296/305)

The table summarizes the statistical properties (i.e., coverage and density) of the scenario boxes describing the vulnerability conditions of each policy regime. The quantitative thresholds defining each scenario box are listed.

substitution and climate sensitivity. Scenario box 1 describes “high climate sensitivity futures,” while Scenario box 2 describes “medium-to-low elasticity of substitution scenarios.” Differences in the vulnerable region exists between these three policy architectures, namely, that the comprehensive GCF policy (“P7. H. Carbon Tax + Co-Tech-R&D[GCF]”) shows a greater area of success than the other three policy architectures.

These results also show that out of the four uncertainties considered in this analysis, (1) elasticity of substitution, (2) climate sensitivity, (3) atmospheric carbon sink capacity, and (4) the discount rate, only the first two determine whether or not the optimal policy achieves the objective of stabilizing CO₂ emissions at sustainable levels before the end of the century. Arguably, out of these two factors, the elasticity of substitution plays a more fundamental role in determining the vulnerability of the policy response, as all scenarios that display medium to low elasticity of substitution are vulnerable across all policy regimes, while high climate sensitivity scenarios induce vulnerability at high elasticity of substitution scenarios for three out of the four policy regimes considered.

On the other hand, the end-of-century 2°C temperature rise target is met in a greater number of futures than the stabilization target. This implies that the former is a more achievable target than the later. Certainly, meeting the stabilization target would be highly beneficial as this would imply that climate change would not be a prevailing public policy problem after the end of the century; however, the results show that this target is met only under very favorable economic and environmental circumstances.

DISCUSSION

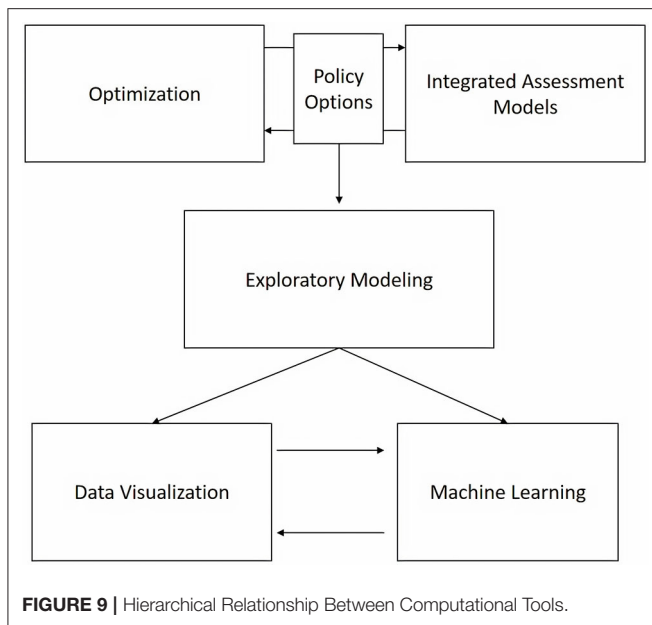
Key Lessons From the Case Study

The results presented in the previous sections show that the combined application of multiple computational intelligence tools produces new insights and more detailed information about

the effectiveness of different climate policy regimes. First, the use of the EDIAM model allows for the joint consideration of multiple regions and the interaction between the economy, the environment, and optimal climate policy. As a result, it is possible to analyze climate change policy multidimensionally in terms of both its ability to mitigate temperature rise and its economic cost (or benefit). Second, by using the EDIAM model in a computational experimentation setting, we show that an uncoordinated carbon tax is the highest cost policy in the majority of cases and that interregional cooperation through the GCF can sometimes be more costly than independent comprehensive climate policy. Our experiment also highlights that there are noticeable differences between policies in terms of the period of time required to achieve stabilization (cooperation between regions generally induces decarbonization faster than non-cooperation). However, we also find that for a considerable number of futures, policy intervention needs to remain in place for as long as 300 years.

Through the application of data visualization techniques, we show that it is possible to describe the dynamics of optimal climate regulation. It doing so, we find that regional differences play a significant role in determining the structure of the optimal policy response. Particularly, we show that in emerging economies carbon taxation is always equal or higher than carbon taxation in advanced economies. In contrast, the technology policy effort of climate policy is stronger in advanced economies than in the emerging economies. Since mitigation technologies are mainly produced in advanced nations, then the optimal policy prioritizes accelerating technology development over taxation in this region, while in the emerging region, higher taxation creates a strong market niche for sustainable energy diffusion, which is used more effectively by R&D and technology subsidies that accelerate the technological catching-up process.

We demonstrate that it is possible to use clustering algorithms to quantitatively identify key drivers of vulnerability of climate policy across various objectives. We find that out of the four



stressors considered, (1) elasticity of substitution, (2) climate sensitivity to GHG emissions, (3) discount rate of economic agents, and (4) carbon sink capacity, only the first two determine whether or not the optimal policy achieves the objective of stabilizing CO₂ emissions at sustainable levels before the end of the century. Considering the relevance of the debate about the appropriate value of the discount rate in climate policy analysis, this finding, which shows that there are more critical drivers of climate policy vulnerability, exemplifies very well the benefits of combining different computational tools for decision analysis in complex systems. Finally, we show that for the independent carbon taxes policy (i.e., P1), in the majority of cases, it is possible to keep the temperature rise below 2°C, but in none of the cases, this policy is able to stabilize CO₂ emissions before the end of the century.

A Hierarchy of Computational Tools for Analyzing Sustainability Challenges

The combined application of various computational tools to this case study yields lessons with respect to their hierarchical relation for analyzing sustainability challenges amid complexity and deep uncertainty. **Figure 9** describes this hierarchy schematically; each block represents an analytical element to be integrated in the analysis of sustainability challenges, and arrows indicate information flows in this hierarchy.

As shown in this case study, the first layer in this hierarchy englobes optimization and Integrated Assessment Models (IAMs). The combination of both perspectives is conducive for analyzing sustainability challenges. IAMs provide the required formalism and tractability for taking into consideration sustainability interdependencies across spheres. Optimization provides the analytical framework needed for formalizing policy options in the light of sustainability objectives. This requires adequate cost estimates of competing alternatives, formalization

of decision restrictions, and sustainability performance metrics for all systems considered. The second layer pertains to the integration of the models produced in the first layer with exploratory modeling (Bankes, 1993; Kwakkel, 2017). The intention of using exploratory modeling is to produce, for each parametrization case, a vector of optimal action. This yields a rich database that maps out changes in optimal action across the often vast ensemble of cases considered. The third layer of this hierarchy connects with the second by the direct application of data visualization and machine learning techniques. Machine learning techniques, in particular clustering techniques and decision rule classifiers, can be used to identify statistically (a) vulnerability conditions of sustainability objectives across policy alternatives and (b) critical thresholds for triggering different actions. Data visualization techniques can be particularly useful to track down changes of the optimal policy response across the parameterization space and to create decision-support tools to be used in participatory planning exercises.

This integration of computational tools is useful because the statistical evidence produced through the integration of these tools leads to a more nuanced understanding of the conditions under which different policy alternatives are more appropriate for achieving sustainability goals. For example, Molina-Perez et al. (2019) apply a similar approach for analyzing sustainability water challenges amid climate, economic, and technological deep uncertainty. In their analysis, the authors integrate econometric, water, and climate modeling tools to develop an IAM, which is combined with an optimization framework that assesses how to best expand the water infrastructure of Monterrey, Mexico. Their results show that it is possible to develop a robust expansion strategy that meets systems' reliability and environmental restrictions without exposing the city to large financial and operational risks. Such strategy is comprised of a diversified collection of projects that considers both conventional and non-conventional expansion strategies and that postpones large infrastructure investment until more information about climate and technological change becomes available.

There are multiple avenues for future research with respect to integrating multiple computational tools for analyzing sustainability challenges. On the one hand, this line research will greatly benefit from standard statistical procedures for designing experimental designs that reduce the risks of biases and increase precision of estimations. This is challenging as each one of these tools (i.e., simulation models, optimization, and machine learning algorithms) needs to be calibrated, trained, and parametrized. In current studies, this is mainly done *ad hoc* and there is little evidence describing, for example, how parameter selection in an optimization routine impacts statistical inference of a classification algorithm; the same is true for experimental designs in exploratory modeling exercises. On the other hand, there is ample room for studying, from a behavioral perspective, how to best transfer findings of these studies to non-specialized audiences. For instance, experimental evidence comparing the impact on knowledge transfer of different combinations of computational

tools could shed light on the most appropriate approach for integration.

CONCLUSIONS

This paper applies DMDU methods to structure an analysis of global climate change mitigation and to demonstrate that the combination of multiple computational tools for analyzing this sort of sustainability challenges leads to richer analytical insights than those produced by traditional monodisciplinary studies.

The scope of the computational experiment in the study considers nine different policy regimes and 600 different optimization cases. The ensemble of cases combines four sources of uncertainty: elasticity of substitution, discount rate, climate sensitivity to GHG, and atmospheric carbon sink capacity. The performance of the different policy regimes is evaluated in terms of the end-of-century conditions. Particularly, the performance of each policy regime is evaluated in terms of its capacity to meet two climate change sustainability objectives: (1) the stabilization of CO₂ emissions and (2) the 2°C temperature rise target.

The analysis shows that the structure of optimal environmental regulation changes markedly across the uncertainty space. The results show that the optimal policy response is most affected by climate sensitivity uncertainty and the elasticity of substitution uncertainty. In particular, the strength of the optimal policy response is directly proportional to the level of climate sensitivity to greenhouse gas emissions and inversely proportional to the elasticity of substitution between the sustainable energy and fossil energy sectors. We also show that the discount rate does affect the structure of the optimal policy response, but its influence is less significant when compared to the influence of climate sensitivity and the elasticity of substitution.

The comparison of GCF-based policy regimes and non-GCF policy regimes shows that the GCF does affect the structure of climate policy. These results show that under the GCF the level of carbon taxation reduces for both regions compared to the level of taxation in the non-cooperative policy regimes. Also under the GCF, the optimal level of effort in R&D and technology subsidies in the emerging region is on average higher than the optimal level of effort in the non-cooperative policy regime. This indicates that under the GCF it is feasible for the emerging region to make higher investments in R&D and technology subsidies and reduce the rate of taxation. Similarly, for the advanced region it is shown that it is possible to reduce the level of carbon taxation by co-funding R&D and technology subsidies in the emerging region.

We use machine learning algorithms to analyze the experimental database. These results show that the objective stabilizing CO₂ emissions below 2°C before the end of the century is rarely met. Two decision relevant clusters describe

this type of vulnerability: (1) high climate sensitivity to greenhouse gas emissions and (2) medium-low elasticity of substitution. In contrast, the 2°C temperature rise target without CO₂ stabilization is met in a greater number of cases. For both types of vulnerability, the role of discount rate in defining the vulnerability conditions is found to be minimal.

This analysis shows that by integrating optimization, complex simulation models, and machine learning algorithms, it is possible to quantitatively identify key drivers of vulnerability of climate change mitigation policies. Drawing on lessons from this case study, we propose an analytical hierarchy of computational tools that can be applied to other sustainability challenges. The first layer of this hierarchy consists of coupling IAMs with optimization to capture sustainability interdependencies across systems and path dependencies of optimal policy decisions. The second layer proposes to use exploratory modeling (Bankes, 1993; Kwakkel, 2017) to deal with deep uncertainty. Finally, the third layer of this hierarchy connects with the second by the direct application of data visualization and machine learning techniques for identifying relevant decision clusters, characterizing vulnerability conditions, and identifying critical sustainability thresholds.

DATA AVAILABILITY STATEMENT

All datasets generated for this study are included in the article/**Supplementary Material**.

AUTHOR CONTRIBUTIONS

EM-P developed the mathematical model used in this study and lead the analysis of experimental results. OE-F developed the computational architecture needed to run the experiment in a cloud computer cluster. HZ-M developed the connection of this study to sustainability sciences and collaborated in the analysis of experimental results. All authors contributed to the article and approved the submitted version.

ACKNOWLEDGMENTS

This paper was based on the doctoral dissertation of EM-P, titled Directed International Technological Change and Climate Policy: New Methods for Identifying Robust Policies Under Conditions of Deep Uncertainty, written as part of the requirements of the doctoral degree in public policy analysis at the Pardee RAND Graduate School.

SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/frobt.2020.00111/full#supplementary-material>

REFERENCES

- Acemoglu, D. (2002). Directed technical change. *Rev. Econ. Stud.* 69, 781–809. doi: 10.1111/1467-937X.00226
- Acemoglu, D., Aghion, P., Bursztyn, L., and Hemous, D. (2012). The environment and directed technical change. *Am. Econ. Rev.* 102, 131–166. doi: 10.1257/aer.102.1.131
- Achtnicht, M., Bühler, G., and Hermeling, C. (2012). The impact of fuel availability on demand for alternative-fuel vehicles. *Transp. Res. Part D* 17, 262–269. doi: 10.1016/j.trd.2011.12.005
- Bankes, S. (1993). Exploratory modeling for policy analysis. *Oper. Res.* 41, 435–449. doi: 10.1287/opre.41.3.435
- Bryant, B. P., and Lempert, R. J. (2010). Thinking inside the box: a participatory, computer-assisted approach to scenario discovery. *Technol. Forecast. Soc. Change*, 77, 34–49. doi: 10.1016/j.techfore.2009.08.002
- Byrd, R. H., Lu, P., Nocedal, J., and Zhu, C. (1995). A limited memory algorithm for bound constrained optimization. *SIAM J. Sci. Comput.* 16, 1190–1208.
- Fischbach, J. R., Johnson, D. R., and Groves, D. G. (2019). Flood damage reduction benefits and costs in Louisiana's 2017 coastal master plan. *Environ. Res. Commun.* 1:111001. doi: 10.1088/2515-7620/ab4b25
- Friedman, J. H., and Fisher, N. I. (1999). Bump hunting in high-dimensional data. *Stat. Comput.* 9, 123–143. doi: 10.1023/A:1008894516817
- Groves, D. G., Bonzanigo, L., Syme, J., Engle, N. L., and Rodriguez Cabanillas, I. (2019b). *Preparing for Future Droughts in Lima, Peru: Enhancing Lima's Drought Management Plan to Meet Future Challenges*. Washington, DC: World Bank. Available online at: <https://openknowledge.worldbank.org/handle/10986/31695>
- Groves, D. G., Kuhn, K., Fischbach, J. R., Johnson, D. R., and Syme, J. (2016). *Analysis to Support Louisiana's Flood Risk and Resilience Program and Application to the National Disaster Resilience Competition*. Santa Monica, CA: RAND Corporation. doi: 10.7249/RR1449
- Groves, D. G., and Lempert, R. J. (2007). A new analytic method for finding policy-relevant scenarios. *Global Environ. Change* 17, 73–85. doi: 10.1016/j.gloenvcha.2006.11.006
- Groves, D. G., Molina-Perez, E., Bloom, E., and Fischbach, J. R. (2019b). "Robust decision making (RDM): application to water planning and climate policy," in *Decision Making under Deep Uncertainty*, eds V. Marchau, W. Walker, P. Bloemen, and S. Popper (Cham: Springer), 135. doi: 10.1007/978-3-030-05252-2_7
- Hull, V., Tuanmu, M.-N., and Liu, J. (2015). Synthesis of human-nature feedbacks. *Ecol. Soc.* 20:17. doi: 10.5751/ES-07404-200317
- IPCC (2013). "Climate change 2013: the physical science basis," in *Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change*, eds T. F. Stocker, D. Qin, G.-K. Plattner, M. Tignor, S. K. Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex, and P. M. Midgley (Cambridge, United Kingdom; New York, NY: Cambridge University Press), 1029–1099.
- Isley, S. C., Lempert, R. J., Popper, S. W., and Vardavas, R. (2015). The effect of near-term policy choices on long-term greenhouse gas transformation pathways. *Glob. Environ. Change* 34, 147–158. doi: 10.1016/j.gloenvcha.2015.06.008
- Kasprzyk, J. R., Nataraj, S., Reed, P. M., and Lempert, R. J. (2013). Many objective robust decision making for complex environmental systems undergoing change. *Environ. Model. Softw.* 42, 55–71. doi: 10.1016/j.envsoft.2012.12.007
- Kwakkel, J. H. (2017). The exploratory modeling workbench: an open source toolkit for exploratory modeling, scenario discovery, and (multi-objective) robust decision making. *Environ. Model. Softw.* 96, 239–250. doi: 10.1016/j.envsoft.2017.06.054
- Lempert, R. J. (2003). *Shaping the Next One Hundred Years: New Methods for Quantitative, Long-Term Policy Analysis*. Santa Monica, CA: Rand Corporation. doi: 10.7249/MR1626
- Lempert, R. J., and Groves, D. G. (2010). Identifying and evaluating robust adaptive policy responses to climate change for water management agencies in the American west. *Technol. Forecast. Soc. Change* 77, 960–974. doi: 10.1016/j.techfore.2010.04.007
- Lempert, R. J., Groves, D. G., Popper, S. W., and Bankes, S. C. (2006). A general, analytic method for generating robust strategies and narrative scenarios. *Manag. Sci.* 52, 514–528. doi: 10.1287/mnsc.10.50.0472
- Liu, J., Hull, V., Batistella, M., DeFries, R., Dietz, T., Fu, F., et al. (2013). Framing sustainability in a telecoupled world. *Ecol. Soc.* 18:26. doi: 10.5751/ES-05873-180226
- Liu, J., Hull, V., Yang, W., Viña, A., Chen, X., Ouyang, Z., et al. (eds.). (2016). "Framing sustainability of coupled human and natural systems," in *Pandas and People: Coupling Human and Natural Systems for Sustainability* (Oxford: Oxford University Press), 15–26. doi: 10.1093/acprof:oso/9780198703549.001.0001
- Marchau, V. A., Walker, W. E., Bloemen, P. J., and Popper, S. W. (2019). *Decision Making Under Deep Uncertainty*. Cham: Springer. doi: 10.1007/978-3-030-05252-2
- Molina-Perez, E. (2016). *Directed International Technological Change and Climate Policy*. Santa Monica, CA: RAND Corporation.
- Molina-Perez, E., Groves, D. G., Popper, S. W., Ramirez, A. I., and Crespo-Elizondo, R. (2019). *Developing a Robust Water Strategy for Monterrey, Mexico*. Santa Monica, CA: RAND Corporation.
- Nordhaus, W. D. (2011). The architecture of climate economics: designing a global agreement on global warming. *Bull. Atom. Sci.* 67, 9–18. doi: 10.1177/0096340210392964
- Ostrom, E. (2009). *Understanding Institutional Diversity*. Princeton; Oxford: Princeton University Press. doi: 10.2307/j.ctt7s7wm
- Ostrom, E. (2011). Background on the institutional analysis and development framework. *Policy Stud. J.* 39, 7–27. doi: 10.1111/j.1541-0072.2010.00394.x
- Papageorgiou, C., Saam, M., and Schulte, P. (2013). Elasticity of substitution between clean and dirty energy inputs—a macroeconomic perspective. *ZEW-Centre for European Economic Research Discussion Paper No. 13–087*. doi: 10.2139/ssrn.2349534
- Popper, S. W., Berrebi, C., Griffin, J., Crane, K., Light, T., and Daehner, E. M. (2009). *Natural Gas and Israel's Energy Future*. Santa Monica, CA: RAND Corporation. doi: 10.7249/RB9476-1
- Pottier, A., Hourcade, J.-C., and Espagne, E. (2014). Modelling the redirection of technical change: the pitfalls of incorporeal visions of the economy. *Energy Econ.* 42, 213–218. doi: 10.1016/j.eneco.2013.12.003
- Taylor, K. E., Stouffer, R. J., and Meehl, G. A. (2012). An overview of CMIP5 and the experiment design. *Bull. Am. Meteorol. Soc.* 93, 485–498. doi: 10.1175/BAMS-D-11-00094.1
- Train Kenneth, E. (2003). *Discrete Choice Methods With Simulation*. Cambridge: Cambridge University Press. doi: 10.1017/CBO9780511753930

Conflict of Interest: The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2020 Molina-Perez, Esquivel-Flores and Zamora-Maldonado. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.



Personogenesis Through Imitating Human Behavior in a Humanoid Robot “Alter3”

Atsushi Masumori^{1,2*}, Norihiro Maruyama^{1,2†} and Takashi Ikegami^{1,2}

¹ Department of General Systems Science, University of Tokyo, Tokyo, Japan, ² Alternative Machine Inc., Tokyo, Japan

OPEN ACCESS

Edited by:

Georg Martius,
Max Planck Institute for Intelligent
Systems, Germany

Reviewed by:

Jun Tani,
Okinawa Institute of Science and
Technology Graduate University,
Japan

Alan Frank Thomas Winfield,
University of the West of England,
United Kingdom

*Correspondence:

Atsushi Masumori
masumori@sacral.c.u-tokyo.ac.jp

†These authors have contributed
equally to this work

Specialty section:

This article was submitted to
Computational Intelligence in
Robotics,
a section of the journal
Frontiers in Robotics and AI

Received: 04 February 2020

Accepted: 09 October 2020

Published: 18 January 2021

Citation:

Masumori A, Maruyama N and
Ikegami T (2021) Personogenesis
Through Imitating Human Behavior in
a Humanoid Robot “Alter3”.
Front. Robot. AI 7:532375.
doi: 10.3389/frobt.2020.532375

In this study, we report the investigations conducted on the mimetic behavior of a new humanoid robot called Alter3. Alter3 autonomously imitates the motions of a person in front of it and stores the motion sequences in its memory. Alter3 also uses a self-simulator to simulate its own motions before executing them and generates a self-image. If the visual perception (of a person's motion being imitated) and the imitating self-image differ significantly, Alter3 retrieves a motion sequence closer to the target motion from its memory and executes it. We investigate how this mimetic behavior develops interacting with human, by analyzing memory dynamics and information flow between Alter3 and a interacting person. One important observation from this study is that when Alter3 fails to imitate a person's motion, the person tend to imitate Alter3 instead. This tendency is quantified by the alternation of the direction of information flow. This spontaneous role-switching behavior between a human and Alter3 is a way to initiate personality formation (i.e., personogenesis) in Alter3.

Keywords: personogenesis, agency, imitation, self-simulation, memory, reconsolidation, humanoid robot

1. INTRODUCTION

We present a new humanoid robot named Alter3 (**Figure 2**) and analyze the dynamics of Alter3's interactions with humans. The philosophy behind Alter3 is grounded in long-running discussions around human/robot cognition (see section 2). We are particularly interested in Rössler's argument of an artificial cognitive map system (Rössler, 1981), and we attempt to realize and extend his ideas with Alter3. Rössler named the self-organization of a dynamic cognitive map under locomotion as the “Helmholtz–Poincaré–Tolman” hypothesis based on Helmholtz's internal map system generated through locomotion (Von Helmholtz, 1867), Poincaré's internal and external representation of the world (Poincaré, 1905), and Tolman, O'Keefe, and Nadal's ideas of a cognitive map, which was later discussed in relation to placing cells in the hippocampus (O'Keefe and Nadel, 1978).

Dayan et al. (1995) later argued that Helmholtz's idea could be implemented in a self-supervised hierarchical neural system, which they called a Helmholtz machine. The Helmholtz machine is based on an inference system that uses variational Bayesian networks. It is essentially equivalent to a Boltzmann machine (Hinton and Sejnowski, 1983) and provides a basis for a variational autoencoder (Kingma and Welling, 2014).

Apart from the probabilistic approach to cognitive map systems, a dynamic systems approach has also been studied. Jun Tani, for example, studied the self-organization of a neural representation of an environment, with a recurrent neural network on a navigation robot in a real environment (Tani, 1996). More recently, using long short-term memory networks, Noguchi et al. (2019) demonstrated the modality of self-organization of a cognitive map in a navigation robot. The

current research is not a probabilistic approach to cognitive map systems. However, it is not, in the strict sense, a dynamic systems approach, as the updates of the entire system are not synchronized, and above all, it can only operate as a system when it interacts with humans.

Rössler's autonomous navigation system is based on a digital scanner and a digital flight simulator. Alter3 is the realization of another autonomous machine, with a completely new purpose. The purpose is to investigate the ways in which a humanoid robot becomes a person, which we call the "personogenesis" (Rossler et al., 2019) of a humanoid robot. "Personogenesis" refers to the process by which an agent acquires free will to act out of its own volition, much like an independent person. In addition, it may perceive happiness from the emotions of a person or be able to display similar emotions. For example, human babies imitate the mother's facial expressions automatically, which is called primitive mimicry (Meltzoff and Moore, 1989), and then advance to the personogenesis phase. In Rossler et al. (2019) and Rossler (1987), this advancement is initiated by two coupled agents: "the two mirror-competent brain equation carriers with cognition and memory and mirror competence suddenly become, if coupled in a cross-caring fashion, their own masters." In other words, coupled agents (one of the two can be a real person) can suddenly share and exchange happy mental states with each other. Our primary goal is to observe the transition from the primitive automatic mimicry phase to personogenesis in a humanoid robot.

Alter3 autonomously imitates the motion of a person in front of it and stores those motions in its memory in the form of a time series. At the same time, the self-simulator included in Alter3 simulates Alter3's motions and generates a self-image. If the visual perception (the motion of the person being imitated) and the self-image differ significantly, Alter3 retrieves a motion from memory that is closer to the human motion and enacts the retrieved motion. In both the cases, Alter3's spontaneous neural dynamics affect the generation of motion. Thus, Alter3 involves three primary functions/features: an automatic mimicry capacity, self-simulation, and memory selection/variation with a neural noise source. To the best of our knowledge, this is one of the first study to focus on memory-driven imitation in a humanoid robot.

1.1. Automatic Mimicry Capacity

Piaget's major assumption in his cognitive development theory (Piaget, 1966) is based on mimicry. It is known that newborn infants automatically imitate the facial and manual gestures of adults (Meltzoff and Moore, 1989). This ability is believed to be an innate characteristic and is observed in human babies when they are approximately 3 months old. In the design of Alter3, imitation is considered an important step in the development of cognitive abilities. Therefore, we implemented an algorithm that imitates the motion of a person captured by the eye camera.

1.2. Self-Simulation

A self-simulator forms a mental image of the self. Recently, David Ha and Jürgen Schmidhuber worked on model-based reinforcement learning and proposed a "world model" (Ha and Schmidhuber, 2018). In this model, an agent learns an

environmental model that includes its behavior and uses the environmental model for simulation. It demonstrates that a control policy can be trained in the simulated world.

While these are examples of self-simulators that include not only the self but also the environment, Alter3's self-simulators are more specific to the self-image. A more pertinent study is that of the self-modeling agent proposed by Bongard et al. (2006). Because a four-legged agent acquires a self-model by autonomously generating its own behavior, even if one of the legs is removed, the self-model is able to adapt. Kwiatkowski and Lipson (2019) extended this study by replacing the self-model with a neural network.

In these studies, the self-simulator is autonomously acquired through evolutionary processes or through learning by neural networks; however, in our study, we assume that the self-simulator has already been acquired in Alter3, and the parameters are fixed. This is done to focus specifically on the acquisition of individuality, based on the development of memory through the imitation of human motion.

1.3. Memory Selection and Variation

As soon as Alter3 generates a motion, it stores the motion pattern in its memory buffer. The memory is realized as a queue of chunks (3 s each), with a size of 50 chunks (= 1,500 frames). When the memory is full, the oldest memory chunks are removed, and new memory chunks are added to the queue (i.e., first in, first out).

Alter3 imitates the behavior of the person in front of it (this is called the awake or open-eye mode). Alter3 uses the memory queue when it is difficult to imitate behavior or when no human is in front of it. It searches for the optimal behavioral pattern evaluated by the optical flow in the memory chunk. When a memory is retrieved and executed, it is modified by the neural state. This allows the memory to be recalled and rewritten without the presence of a person. Specifically, after the recalled motion is executed, it is combined in spontaneous neural activity to be stored as a slightly different motion. The more it is recalled, the more the memory makes a slightly deformed copy of itself. It can be seen as a Darwinian evolutionary process of the memory. This is called the dream mode or the closed-eye mode.

The details of these algorithms are given in section 3.

2. RELATED WORKS ON IMITATION IN HUMANOID ROBOTS

Imitation of human behavior by humanoid robots is a long-standing theme in terms of cognitive and biological aspects (see e.g., Schaal, 1999). There are two types of imitation studies in robotics: one for learning and the other for communication. Both share the same underlying mechanism of imitation, while the former uses imitation as a learning tool with an explicit purpose, the latter has no specific purpose for imitation besides communication.

Schaal (1999) claimed that imitation would be a promising approach for developing cognition in a humanoid robot. In the recently surveyed article by Hussein et al. (2017), learning

through imitation is presented as a viable research area for novel learning methods. Although most works on imitation consider it as a strategy for learning from humans unidirectionally, we are more interested in bidirectional imitation learning—human to robot and robot to human. We call this approach “imitation for communication.”

Through communication, people develop the social ability to think about others and maintain a good relationship, and imitation plays a significant role in this process. As in Trevarthen’s experiments (Trevarthen, 1977) with infant–mother communication, and Nadal’s study on pretend-play behavior between two children, imitation is a strong driving force for organizing lively interactions (Nadel et al., 2004). Christopher Nehaniv and Kerstin Dautenhahn edited a book on imitation and social learning (Nehaniv and Dautenhahn, 2007). They also started the Aurora project, which aims to help autistic kids acquire social skills with the use of robots (The AuRoRA Project, 1998).

Along with the “imitation for communication” approach, Ikegami and Iizuka (2007) and Iizuka and Ikegami (2004) studied a turn-taking game to show how imitation emerges as a by-product of mutual cooperation. The present work is a continuation of the previous approaches, in a new humanoid body, with new memory dynamics and a self-simulator.

3. SYSTEM ARCHITECTURE

Figure 1 shows an overview of Alter3’s internal system. The system is a combination of Rössler’s autonomous cognitive map system (Rössler, 1981) and Frith, Blakemoore, and Wolpert’s comparator model (Frith et al., 2000). We extended it to include a memory state and a neural network as a spontaneous dynamics circuit. As mentioned earlier, the system is constructed with three functionalities in mind:

- (1) Automatic imitation capability.
- (2) Self-simulation.
- (3) Memory selection and variation through spontaneous dynamics.

In this section, we explain the methods used to achieve the above three functionalities and describe Alter3’s hardware.

3.1. Humanoid Alter3

Alter3’s body has 43 movable air actuator axes, and its motions can be controlled through a remotely placed air compressor that is mediated by a control system (**Figure 2**). More specifically, its motion is controlled by two types of commands: SETAXIS and GETAXIS. A SETAXIS command, which can be regarded as a motor command, is used to set each axis of the humanoid

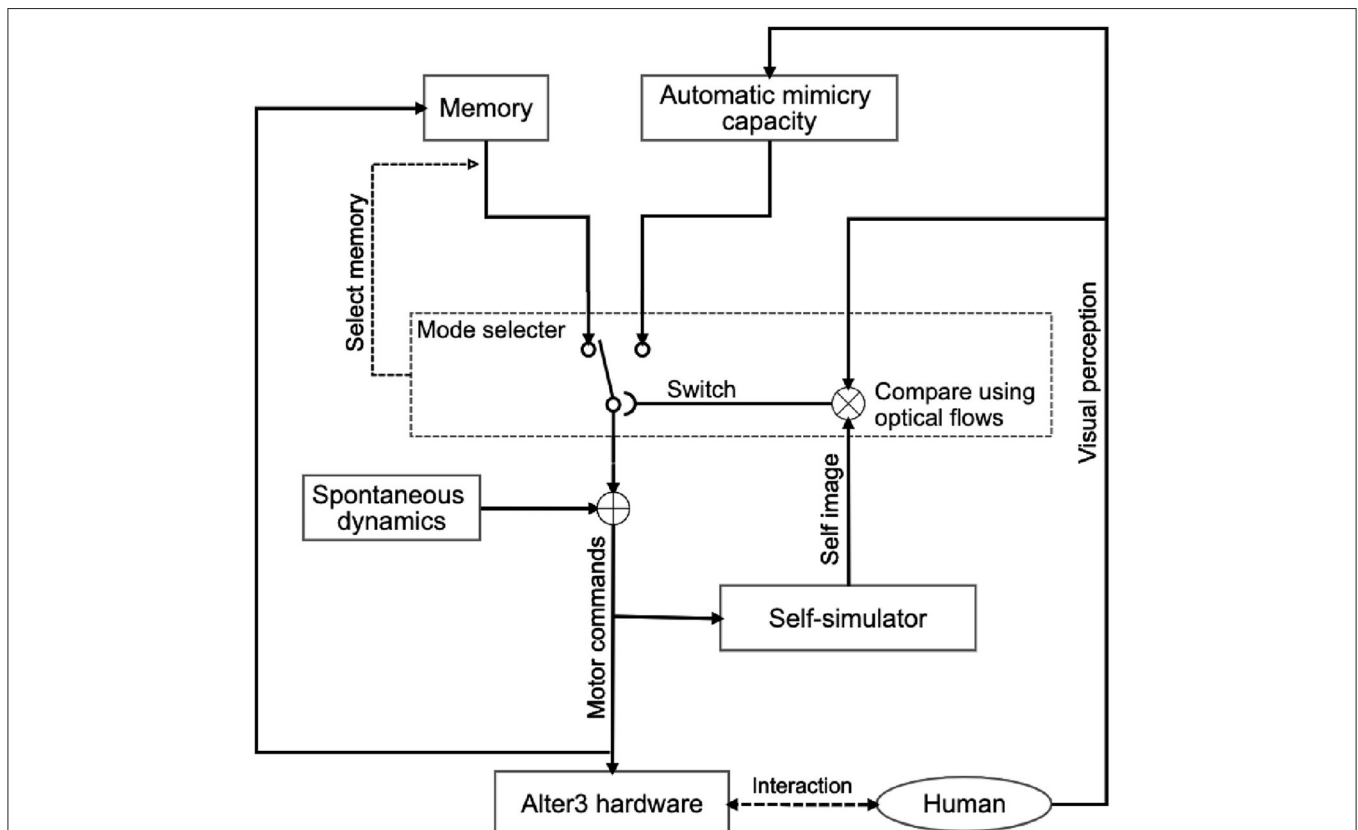


FIGURE 1 | System architecture of Alter3 for the imitation of human behavior. Alter3’s motion is controlled by three main subsystems: a self-simulator, an automatic mimicry unit, and memory storage. Additionally, autonomous neural dynamics perturb the memory system. When Alter3 retrieves a memory chunk and executes it, the retrieved chunk is varied with the neural states and stored again. The details of each module are described in section 3. The mode-selection mechanism is also described in the section and **Figure 5**.

robot to a desired value. By contrast, a GETAXIS command is a command used to retrieve the current axis angle realized on Alter3. Ideally, it is expected that the value obtained from GETAXIS will be the same as the value set by SETAXIS. However, the actual value set for each axis can differ from the intended value. Such differences are caused by physical constraints and latency owing to the body being driven by air actuators. The control system sends commands via a serial port to control the body. Alter3's motions are determined online, and the refresh rate is 100–150 ms.

Alter3 has two cameras, one in each eye, which send visual images to a control system. The camera images are used to extract the key points of the skeleton posture of a human in front of Alter3, using a software called OpenPose (Cao et al., 2017). Alter3 uses the key points of the skeleton to imitate the human posture. In the following sections, the image processing system used for imitation is described in detail.

3.1.1. Automatic Mimicry Capacity

In the awake mode, Alter3's motor commands are generated by the automatic mimicry module through the following processes:

1. Detect a human pose.
2. Map the detected human pose to the angles of the axes.
3. Generate motor commands from the obtained angles and Alter3's spontaneous neural dynamics.

An image from the eye camera is taken as input to a pose detection algorithm. We used OpenPose (Cao et al., 2017) as the algorithm. It detects human poses and generates the positions of key points, such as the head, neck, shoulders, elbows, and wrists. The configuration of the key points of a human skeleton differs from that of the axes in Alter3, and angles of the axes are required as motor commands for Alter3; therefore, we map the positions to the angles. The components responsible for these processes partially constitute Alter3's body schema and can be regarded as the controller in the comparator model (Frith et al., 2000). When

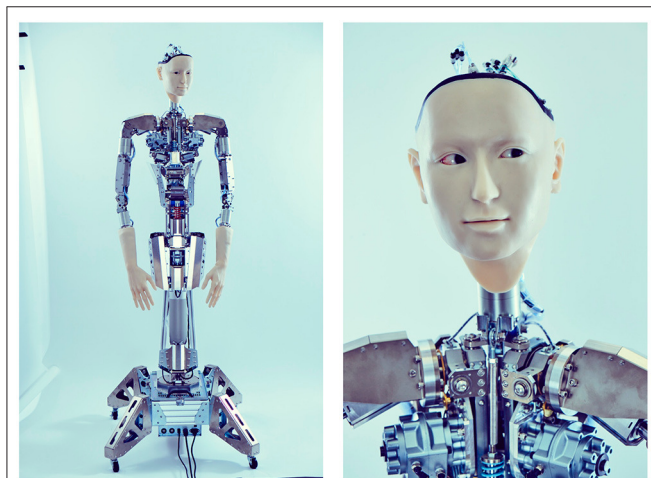


FIGURE 2 | Body of Alter3. The body has 43 axes that are controlled by air actuators. It is equipped with a camera inside each eye.

OpenPose detects poses of multiple people, Alter3 focuses on the center-most person in its visual field and imitates the person's pose. Once the person is locked into Alter3's vision, the person is tracked until the person disappears from its view.

Alter3's spontaneous dynamics consist of spiking neurons (Izhikevich, 2003) that are combined with the calculated angles of the axes as a weighted average to calculate the final axis values (see details in the following sections). The final values are sent to Alter3 as motor commands at every frame, and Alter3 behaves in accordance with the motor commands. Thus, Alter3 not only imitates human motion but also modifies its own motion to an extent based on its spontaneous dynamics.

It should be noted that the choice of whether Alter3 imitates human motion based on the above-mentioned process (awake mode) or based on its memory (dream mode) depends on the result of the comparison between its self-simulation and current visual perception, as described below.

3.1.2. Self-Simulation

Alter3 contains a self-simulator that simulates a future self-image before executing motor commands. The self-simulator is a robot simulator that receives each joint angle as a motor command (which is the same as the SETAXIS command described above) and returns a posture as a visual image (Figure 3). We used a custom-built simulator that visualizes the results of forward kinematics by calculating joint positions from joint angles without a physics engine, other than simple inertia. As Alter3's axes are controlled by air actuators that do not have sufficient torque to control the axes precisely, the actual motions differ from the motor commands. Thus, we manually calibrated the upper and lower limits of the joint angles in the simulation

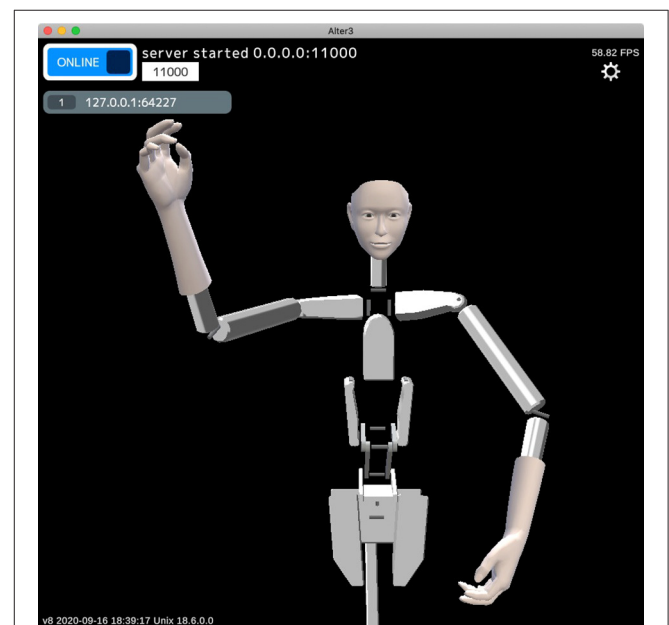


FIGURE 3 | Example of an internal image generated by the self-simulator. The self-simulator receives the SETAXIS commands (motor commands) and generates a visual image.

by comparing the simulated poses and the actual poses by Alter3. This self-simulator can be regarded as a predictor in the comparator model (Frith et al., 2000), which predicts a future state from an efference copy.

The predicted future self-image is compared with the visual perception of the optical flow values. The difference between the two is used to determine the operation mode of Alter3. If the difference between the optical flow values and the predicted self-image exceeds a threshold, the mode switches from awake mode to dream mode, i.e., Alter3 will stop using its automatic mimicry capacity (OpenPose and its mapping function) and will begin using its memory to generate new imitation behavior. The details of this process are explained in the following subsection.

Therefore, Alter3 uses the self-simulator to predict a future posture from the motor commands generated by the automatic mimicry module before executing the commands. It then determines whether it should execute these commands or use memory to imitate the human motion (based on a comparison between the state predicted by the self-simulator and a target human motion).

3.1.3. Memory Selection and Development

Alter3 has a fixed memory size in which the sequence of movements is divided into short chunks that are stored over time. Each memory chunk is a short sequence of behavior but is labeled by an abstract representation of the visual image of the movement. Specifically, we used the optical flow of the self-image for this purpose. When Alter3 identifies that the automatic imitation of a human is not viable under certain criteria, it searches for the optimal movement in its memory by using the labels. In addition, the movement that is retrieved is stored in the memory as a new memory chunk, which allows the formation of a closed loop.

Alter3 stores the executed motor commands in its memory as a memory chunk for every 30 frames. As mentioned in the subsection above, the sequence of motor commands is converted to a self-image via the self-simulator. They are then converted to a series of optical flows. We adapted a dense (lattice) type algorithm to calculate the optical flow. It was originally a two-dimensional vector field, but we adapted it as a scalar field by using the magnitude of the vector. The memory chunk containing 30 frames of the pose sequence was labeled with the time average of the optical flow. Here, we considered the time average of the optical flow as the short-term meaning or label of appearance of the self-motions. For example, when Alter3 performs the action "raising left hand," the motor command is a high-dimensional time series and contains a large amount of information that is irrelevant to the meaning of the motion. It is assumed that the spatial pattern of the optical flow will always take a high value near the upper right side of the body in such cases. Thus, optical flow is qualified as the meaning or the label. In our experiment, optical flow was calculated using the algorithm proposed by Farneback (2003), and OpenCV library (Bradski, 2000) was used for the actual implementation. The memory was realized as a queue of memory chunks, and its size was limited (50 chunks = 1,500 frames). Thus, if the memory was full, the oldest memory chunk was removed, and a new memory chunk was added (i.e., first in, first out). **Figure 4** shows this process.

Alter3 can replay past motions based on memory in the dream mode. This memory recall and motion replay occurs in the following two cases.

1. When no human is in sight.
2. When a self-simulated motion differs significantly from the target human motion.

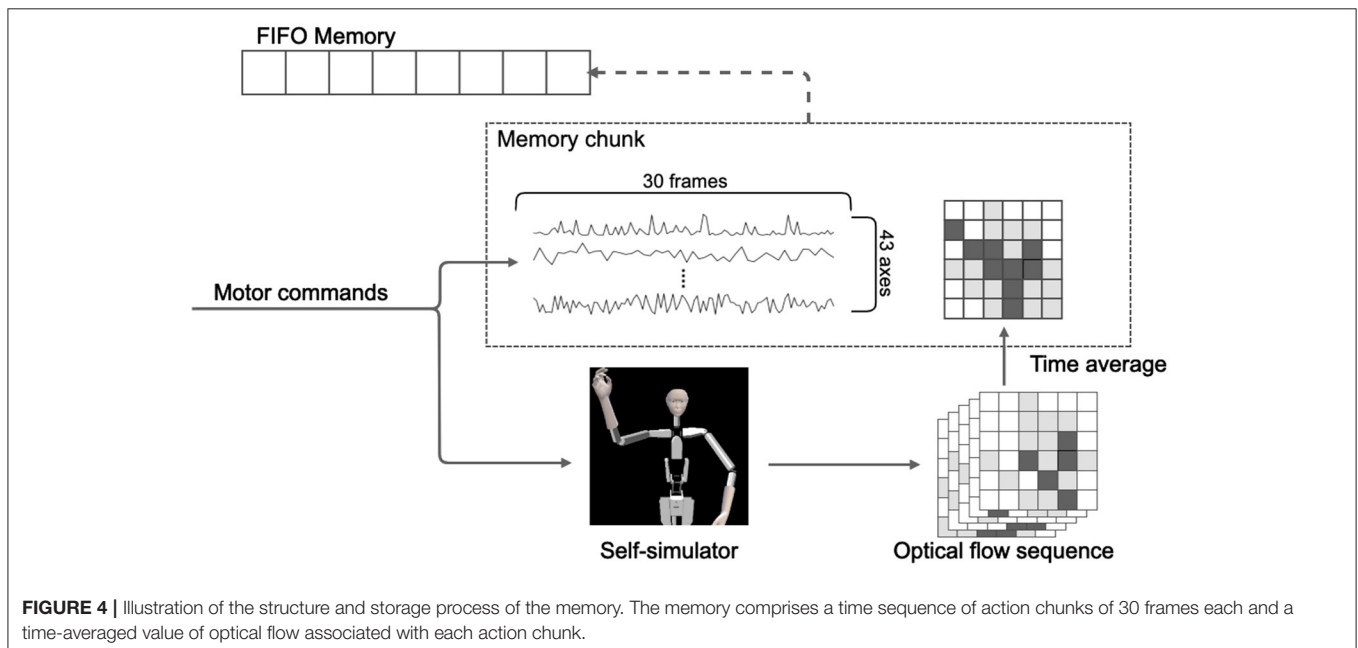


FIGURE 4 | Illustration of the structure and storage process of the memory. The memory comprises a time sequence of action chunks of 30 frames each and a time-averaged value of optical flow associated with each action chunk.

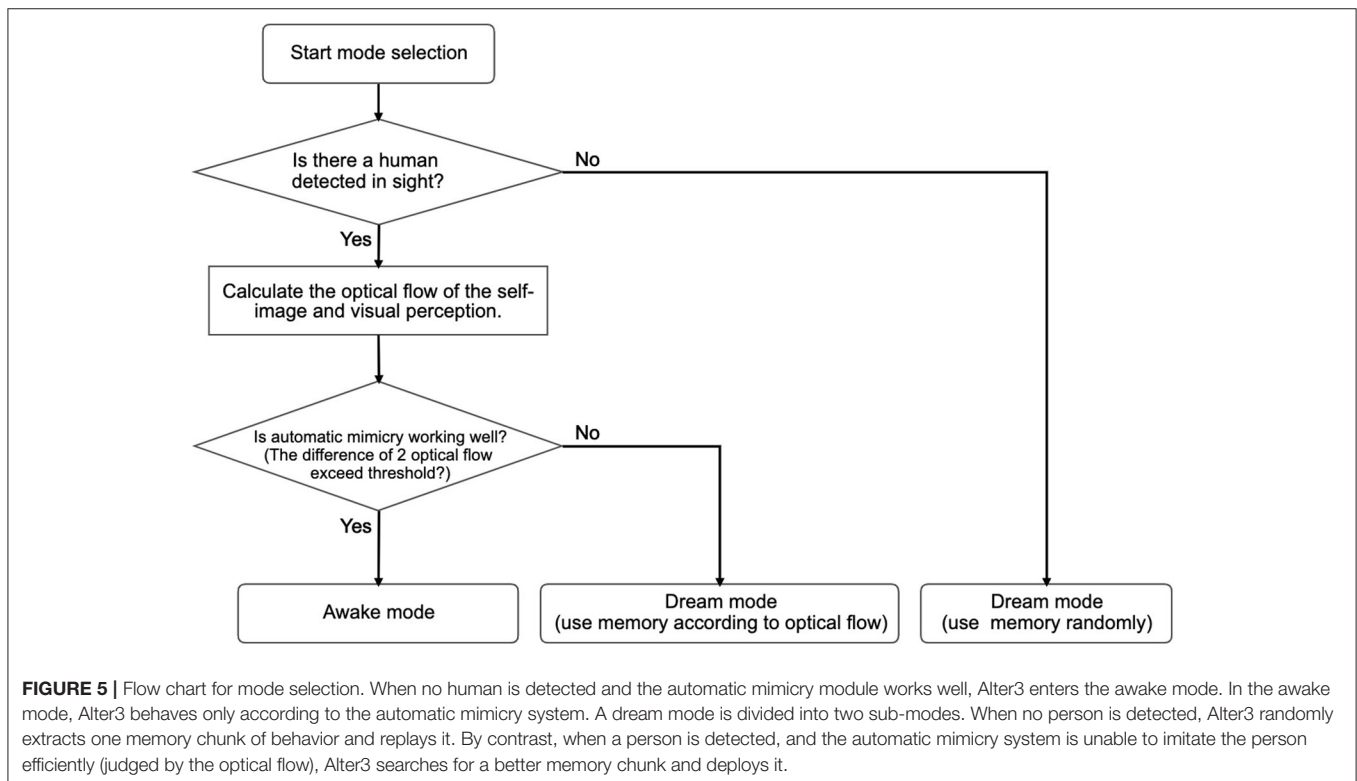
The first case is specifically defined for when OpenPose detects no humans for 100 frames. In this case, Alter3 recalls a motion sequence randomly from memory and replays it. When replaying the motion, Alter3's spontaneous dynamics, which consist of spiking neurons, causes a minor change in the motion as with the case of automatic mimicry (the details of this mutation process are described in the next subsection). The mutated motion is then stored as a new memory. In this case, the memory is reconstructed by store-replay cycles and spontaneous dynamics, without any inputs from the environment. This is similar to memory consolidation in a dream, where memory is reactivated and reorganized (e.g., Wamsley et al., 2010). When a human comes into sight, Alter3 switches to the awake mode.

The second case is specifically defined for when the difference between the optical flows of the self-simulated visual images and the optical flows of the visual perception (human image) exceeds a certain threshold during a short period (15 frames). Mean squared error is used to measure the difference between the two optical flows. In this case, Alter3 retrieves a memory chunk that has been labeled with the optical flow values that are closer to those of the current visual image from the camera and replays the motion. The replayed motion is also mutated by the spontaneous dynamics. The motion is labeled as having a certain optical flow and is stored as a new memory. This is similar to memory reconsolidation, where the recalled memory becomes temporally unstable; then, the memory is consolidated again and becomes stable (e.g., Suzuki et al., 2004). If the optical flow values of the recalled motion are close to the values of the current human motion when a memory chunk is replayed, then Alter3 switches back to awake mode. The algorithms for mode selection are summarized in **Figure 5**.

It should be noted that both memory recall mechanisms explained above are not simple replay mechanisms. Rather, both are memory reconstructions with mutations that are caused by spontaneous dynamics. We expect that the memory recall mechanisms will allow Alter3 to explore new movement patterns that cannot be generated from its automatic mimicry capacity. Additionally, the second recall mechanism can select memories in accordance with the ability to imitate humans, for a given memory chunk; therefore, it develops the contents of memory according to the imitation ability. As a result, we expect that memory structures can evolve through the experimental imitations of human agents.

3.2. Memory Variation by Spontaneous Dynamics

Alter3 has internal spontaneous dynamics that act as a central pattern generator (CPG). This generator has no input from the environment. It consists of spiking neurons (see **Appendix** for the details of the neuron model). The first reason for using spiking neurons instead of other chaotic dynamical systems or stochastic dynamic systems is that we intend to add a learning process with stimulus input in the future work (e.g., the difference between simulated future self-image and target human motion might be used as stimulus input to the spiking neurons). The second reason is that, in this research, it is important that memory becomes unstable with the internal dynamics when it is recalled, i.e., the dynamics are used to perturb the memory. Thus, it would be better if the dynamics kept changing with synaptic plasticity. We compared the dynamics of spiking neurons with synaptic plasticity to spiking neurons without synaptic plasticity and random patterns. The results (**Figure A1**) show that the



generated patterns of the spiking neurons with synaptic plasticity were more structured and temporally richer than the ones without plasticity (see **Appendix 2** for the details of this analysis). For these reasons, we adopted spiking neurons with plasticity as the candidates for noise sources to perturb memory.

The dynamics of the CPG are added to the motor commands before they are sent to Alter3, which implies that the dynamics also mutate recalled memories, much like memory reconsolidation. The original motor commands generated by automatic imitation or memory selection are always affected by the CPG. Specifically, final motor commands realized by Alter3's hardware are taken in a weighted summation of the original motor commands and output of the CPG. We set the weight of the CPG output to 0.1, and the weight of original motor commands to 0.9. In other words, CPG dynamics mutate recalled memories, like memory reconsolidation.

4. EXPERIMENTS

We conducted experiments with Alter3 at the NRW-Forum, Düsseldorf between April 26 and May 4, 2019. During the

experiments, Alter3 was located in the exhibition room (**Figure 6**, left), which is a public space. The public could freely visit the exhibition and witness Alter3's movements. They were allowed to interact with it through their own movements (**Figure 6**, right; see also **Supplementary Video 1**). There was no limitation on the duration for which a person can interact with Alter3, and no information about the experiment was provided besides the fact that Alter3 could imitate human motion. The advantage of a public demonstration was that people of all ages, genders, and nationalities could come to see Alter3. Furthermore, as our policy was to experiment with robots in an open and natural environment, the demonstration was a welcome activity. It is also possible to conduct longer experiments, which can last for weeks (Ikegami, 2010, 2013; Masumori et al., 2020).

We performed six experiments, each consisting of 100,000 frames and lasting approximately 4–5 h. During the experiments, we recorded Alter3's motor commands, its actual motion data, and the human motion data (**Figure 7**). We analyzed these data to understand how Alter3's behavior changed during the experiments.

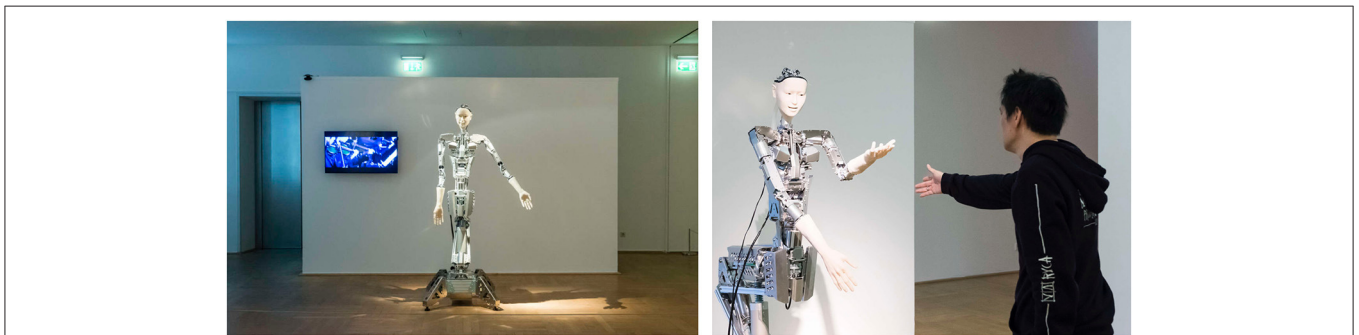


FIGURE 6 | Alter3 at the exhibition NRW-Forum, Düsseldorf. It was evident to the public that Alter3 was trying to imitate the pose of a person.

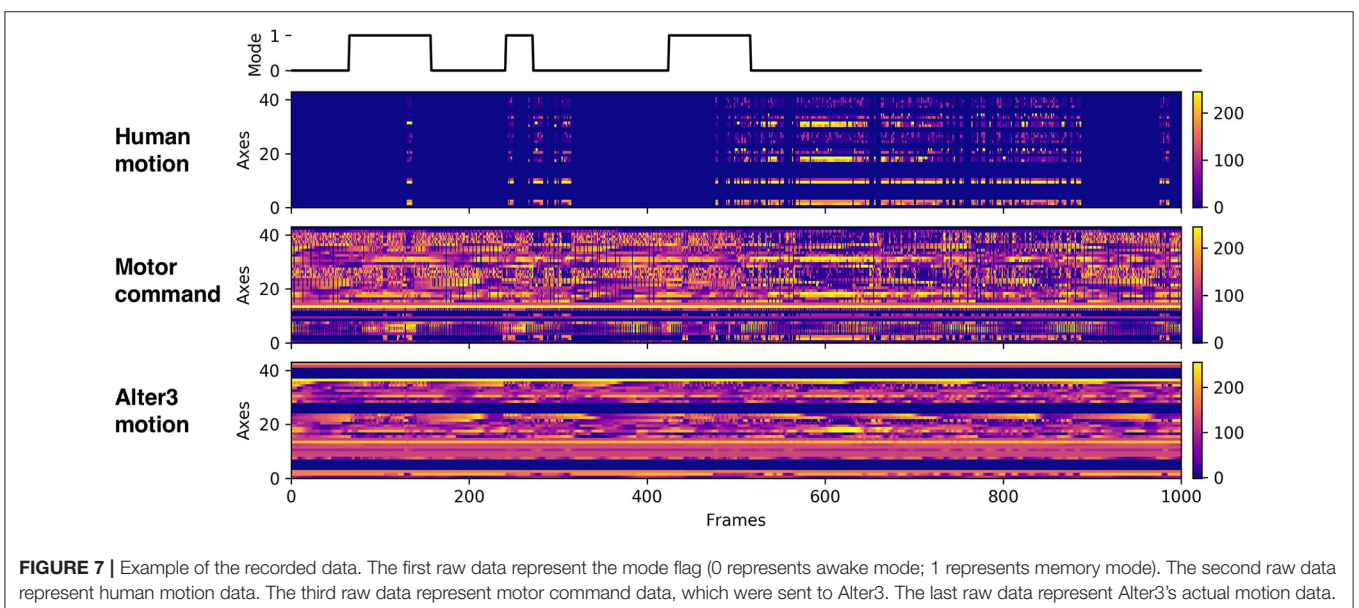


FIGURE 7 | Example of the recorded data. The first row data represent the mode flag (0 represents awake mode; 1 represents memory mode). The second row data represent human motion data. The third row data represent motor command data, which were sent to Alter3. The last row data represent Alter3's actual motion data.

5. RESULTS

5.1. Development of Memory Structure

We analyzed the change in memory and actual motions of Alter3. The memory and actual motion values (values of SETAXIS and

GETAXIS) have 43 dimensions; hence, we adapted a dimension-reduction algorithm called UMAP (McInnes et al., 2018) to visualize them. **Figure 8** shows the results of the dimension-reduction by UMAP, which reduced the memory and actual motion data of Alter3 to two dimensions. These results show

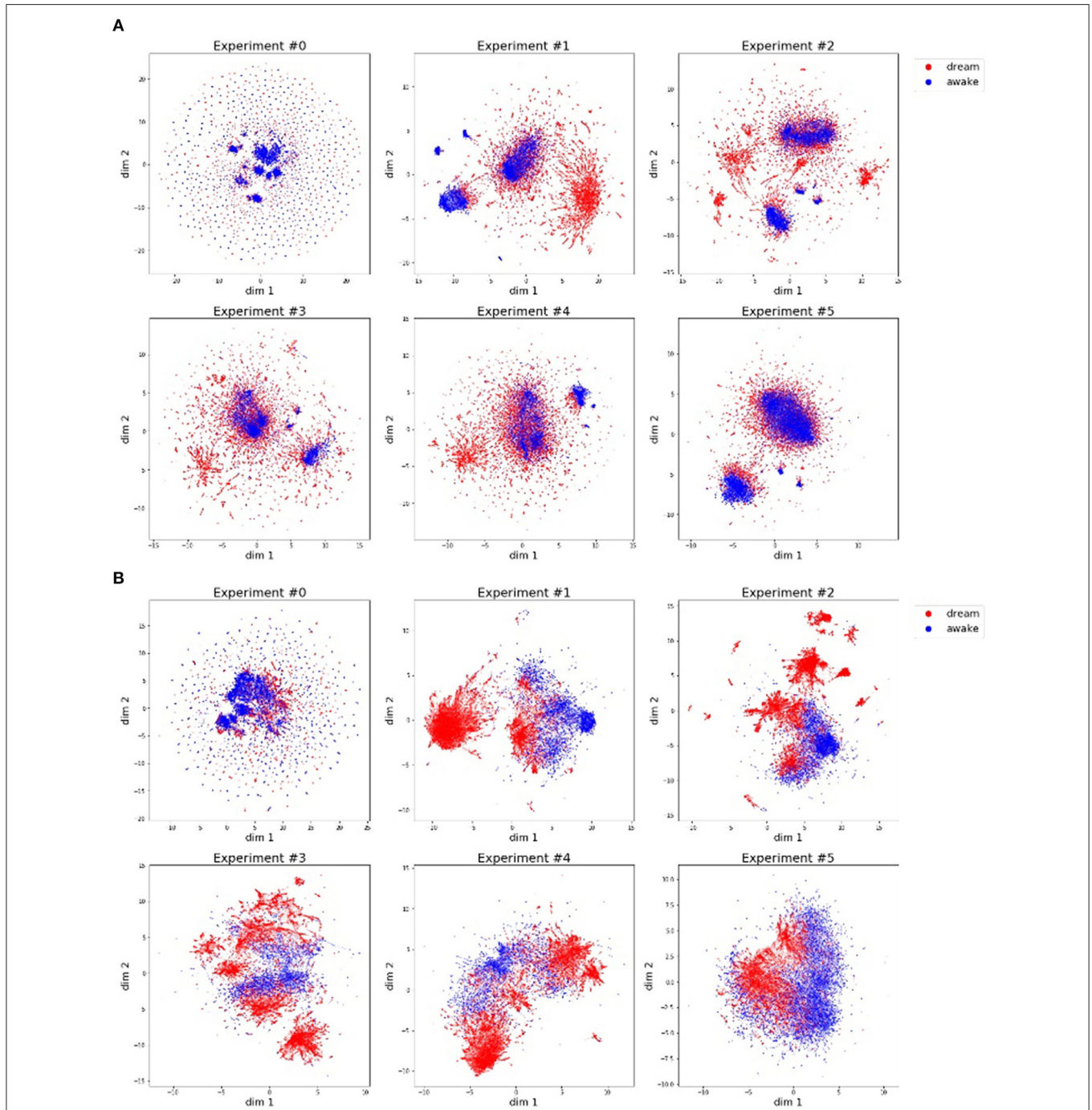


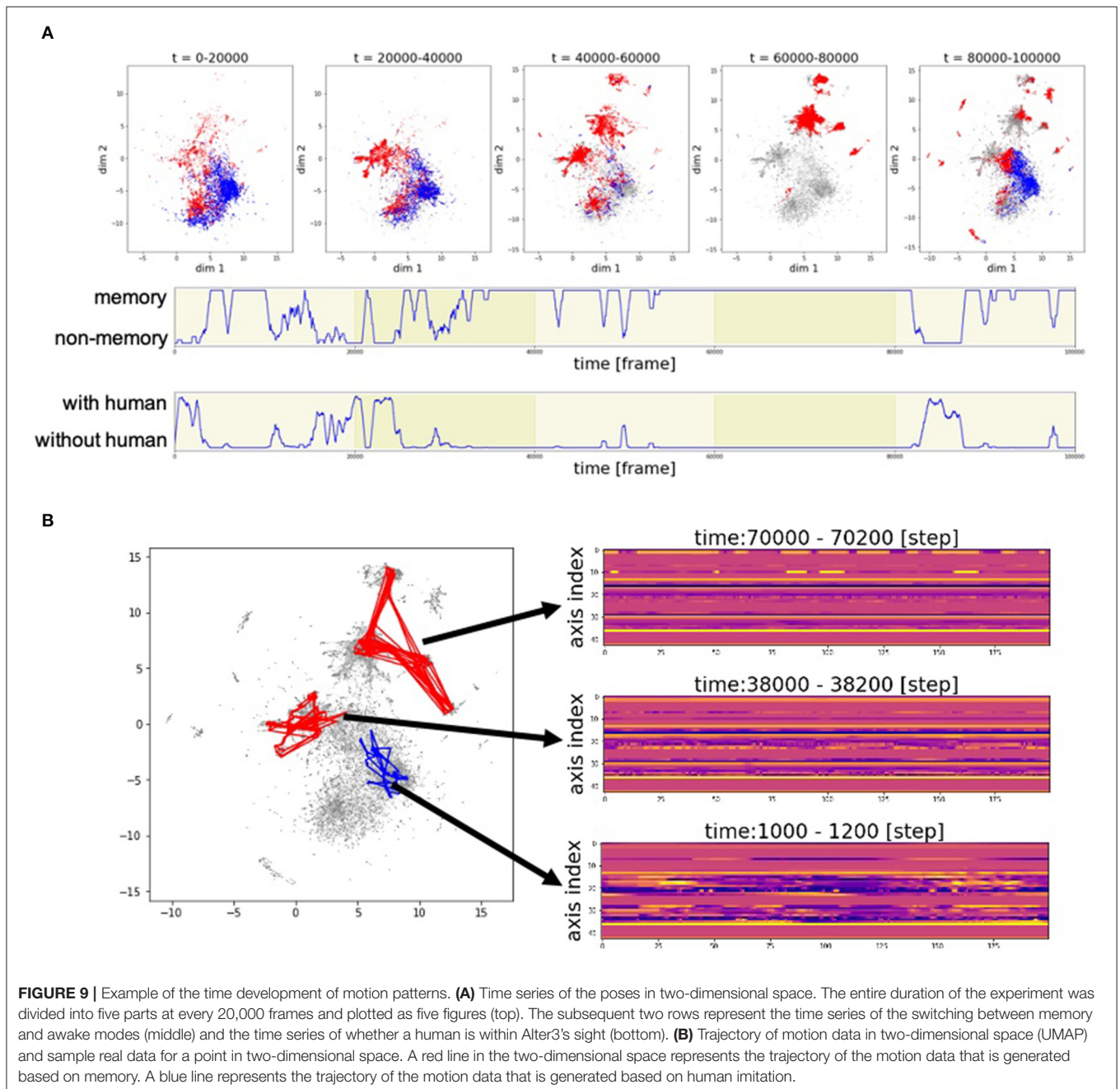
FIGURE 8 | Motor commands **(A)** and motion data **(B)** are projected onto a two-dimensional space using the dimension reducing algorithm, UMAP. The blue dots indicate that the pose is generated by imitating human motion, and the red dots represent poses generated from memory. **(A)** Motor commands data (history of Alter3's motor commands: SETAXIS) for each experiment. **(B)** Motion data (history of Alter3's actual motion: GETAXIS) for each experiment.

a different pattern for each experiment, especially experiments #0 and #5. We consider that these differences reflect differences in the interactions between Alter3 and humans. This suggests that different personalities in Alter3 emerge from different environments (e.g., differences in the frequency and duration of people's stays and motion patterns).

As shown in this figure, in almost all the experiments, the poses generated in the awake mode and the poses generated in the dream mode have different clusters, although some of these clusters overlap. The former poses tend to have more clusters than the latter ones. This suggests that Alter3 not

only copied human motions but also varied them using its memory mutation and selection process. The memory data (Figure 8A) and the actual motion data (Figure 8B) reflect the same tendencies. However, they also marginally differ because of Alter3's construction: Alter3's axes are controlled by air actuators, and they do not have sufficient torque to control the axes precisely. Thus, the actual motions differ from the motor commands.

Figure 9 shows the developments in the motion patterns over time. Figure 9A (top) shows that the clusters of the poses generated from memory, represented by the red dots, are initially



located near the clusters of the poses generated by automatic mimicry capacity, represented by the blue dots. Then, the red clusters begin to vary and move away from the blue clusters. At 40,000–60,000 frames and 80,000–100,000 frames, many red clusters can be observed. In these phases, there are cases where Alter3 retrieves a memory and behaves accordingly despite a person being in its sight (Figure 9A, middle and bottom). In such a case, memory selection and the reconsolidation process occur. These results suggest that the memory selection and variation process work well to diversify memory, rather than just copy human motion.

The motion pattern of Alter3 can be represented in a two-dimensional plane. Figure 9B shows the trajectories of the motion data in two-dimensional space (UMAP), and some samples of the data of actual points in the two-dimensional space. It can be observed that the complex motion pattern derived from human motion (at 1000–1200 frames) gradually converges to relatively static motions (at 38,000–38,200 frames and at 70,000–70,200 frames), probably because there were few humans in Alter3's sight at 38,000–38,200 frames, and none at 70,000–70,200 frames. This suggests that Alter3's memory diversifies itself through interactions with the environment (human) at first. However, without such interactions, its memory is overwritten by its spontaneous activity and gradually disappears, similar to forgetting dynamics in actual humans.

5.2. Information Flow Between Alter3 and Human

To evaluate whether Alter3 could effectively imitate human motion and whether humans also imitated Alter3, we analyzed the information flow between Alter3 and humans. We used transfer entropy (TE) to estimate the information flow between the motions of Alter3 and the humans during the experiments. TE measures directed information transfer (Schreiber, 2000). A high TE from one entity to another indicates that the former affects the latter. Thus, TE enables us to estimate causation during an imitation.

The TE from time series J to time series I is defined as

$$TE_{J,I} = \sum p(i_{t+1}, i_t^{(k)}, j_{t+1}^{(l)}) \log \frac{p(i_{t+1} | i_t^{(k)}, j_{t+1}^{(l)})}{p(i_{t+1} | i_t^{(k)})}, \quad (1)$$

where i_t denotes the value of I at time t , j_t denotes the value of J at time t , and i_{t+1} denotes the value of i at time $t + 1$. Parameters k and l give the order of the TE and represent the number of time bins in the past that are used to calculate the histories of time series i and j . Here, they are set to $k = l$ and $k = 3$.

We computed the TE between the motion data of both Alter3 and humans (continuous multivariate data) using the Kraskov–Stögbauer–Grassberger estimator in the JIDT library (Lizier, 2014) and compared the results for the awake and memory conditions. The awake condition was defined to be equivalent to the awake mode explained above. The memory conditions were defined such that there was a human in front of Alter3, but the error of the optical flow exceeded the threshold, and memory was used to generate Alter3's motion.

The mean TE values between Alter3's motion and human motion are shown in Figure 10. In the awake mode, the value of TE from Alter3's motion to human motion was significantly lower than in the opposite direction (Mann–Whitney U -test, $n = 6$, $p = 0.0025$). This implies that information flow from humans to Alter3 was higher than the flow from Alter3 to humans. This suggests that Alter3 could imitate human motion effectively. In contrast, for the memory condition, the value of TE from Alter3 to human motion was significantly higher than the TE value for the opposite direction (Mann–Whitney U -test, $n = 6$, $p = 0.0227$). This suggests that information flow was reversed in the memory condition, and humans tended to imitate Alter3. During the dream mode, the motions were selected from memory based on the similarity of the visual image-based motion pattern (optical flow) between the poses of Alter3 and a human, rather than the similarities of joint angles itself. Thus, under this condition, the similarity of the motion at the joint angle level will not necessarily be as high as it would be in the awake mode. Such a difference may induce people to start imitating Alter3.

TE varies temporally. As an example of a time series, Figure 11 shows an alternation of local TE between Alter3 and human motions. It shows that, in the memory conditions, the local TE from Alter3's motion to human motion was often higher than in the opposite direction. In addition, in the awake mode, the local TE from Alter3's motion to human motion was sometimes higher than that from human to Alter3. These results imply that the causes and effects of the imitation were often reversed over time; thus, Alter3 and humans imitated each other. We think that this

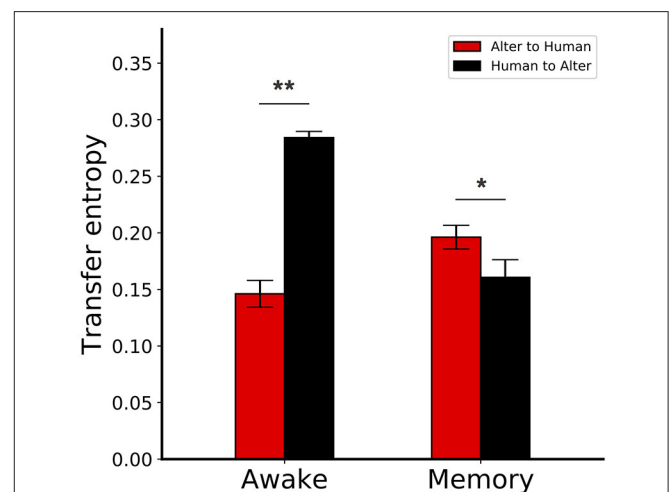


FIGURE 10 | Transfer entropy (TE) between Alter3 and human motions. In the awake condition, the TE from Alter3's motion to human motion was significantly lower than the reverse case. In contrast, for the memory condition, the TE from Alter3 to human motion was significantly higher than the TE in the opposite direction. The awake conditions were equivalent to the awake mode, where Alter3 imitated human motion with its automatic mimicry module. The memory conditions were defined when Alter3 used memory to generate motion (i.e., there was a human in front of Alter3, but the difference between the optical flow values of the human motion image and the simulated future self-image of Alter3 exceeded a threshold; thus, memory was used to generate motion). * $p < 0.05$, ** $p < 0.01$.

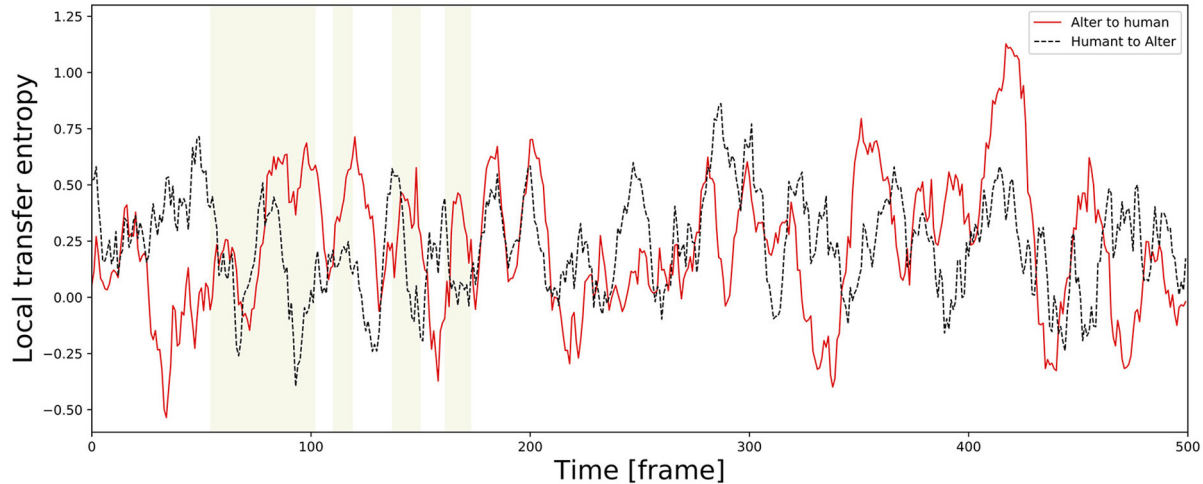


FIGURE 11 | Example of a time series of local transfer entropy (TE) between Alter3 and human motions. The yellow zones indicate a memory condition in the dream mode, where Alter3 used its memory to generate a posture. In the memory condition, the values of local TE from Alter3's motion to human motion tend to be higher than those from human motion to Alter3's motion. Furthermore, in the awake mode, in which Alter3 imitates human motion by automatic mimicry capacity, there were also cases where the values of local TE from Alter3 to human motion were sometimes considerably higher (e.g., close to 420 frames in the figure). We think that this was because Alter3's generated motion pattern was sometimes not as good as expected, thus Alter3 failed to imitate human motion. It seems that this situation led people to imitate Alter3 in turn, thus TE values from Alter3 to human sometimes higher than opposite direction even in the awake mode.

was because Alter3's generated motion pattern was sometimes not as good as expected, thus Alter3 failed to imitate human motion. It seems that such a situation led people to imitate Alter3 in turn, thus the TE value from Alter3 to human was sometimes higher than opposite direction.

6. DISCUSSIONS

Alter3 is programmed to imitate the motion of a person in front of it. A human pose detection algorithm (OpenPose) extracts the key points of the skeleton from the posture pattern. However, Alter3 sometimes fails to imitate the motion. The imitation rating is based on the difference between the optical flow pattern of Alter3 and the optical flow pattern of the person Alter3 attempts to imitate. The smaller the difference, the better the imitation. The main reasons why Alter3 sometimes fails to imitate human motion are (i) physical constraints imposed by the mechanical structure of Alter3, (ii) incorrect detection caused by OpenPose or disturbance to the eye camera, (iii) the dynamic characteristics of Alter3's unstable process, (iv) a significant time delay between the control program and the motor output, and (v) Alter3 encountering a style of motion that cannot be imitated. Such types of failures play an important role for Alter3, such as organizing memory through selection and mutation processes and inducing role switching in interactions with human, as discussed below.

Introducing memory into Alter3, we incorporated an imitation recovery process: if Alter3 fails to imitate human motion with automatic mimicry capacity, Alter3 uses memory to imitate the motion. Alter3's spontaneous neural dynamics commonly affects the generation of motion. Therefore, posture

patterns are not only stored in memory but also changed over time. Owing to the selection and the mutation processes, memories are copied and changed when they are used. If Alter3 uses a stored pattern frequently, more copies of this pattern emerge with modifications.

Alter3's organized motion is generated by the automatic mimicry capacity or through Alter3's memory. Therefore, the whole posture space of Alter3 is decomposed into two categories. One consists of the postures provided by estimating human postures, and the other category has the self-organized postures generated through memory selection and variation. The decomposition is shown by applying the UMAP compression in **Figure 8**. These two categories are created spontaneously through interaction with humans. Moreover, if no person appears in front of Alter3 for a certain period, the postures in the latter category gradually change and converge to Alter3's spontaneous dynamics provided by the spiking neurons, after which another category is organized.

To determine whether Alter3 imitates people's postures or whether people imitate Alter3, we measured the TE between Alter3 and the people whose motions it seemed to imitate. The results suggest that people often imitate Alter3 strongly when Alter3 is in the dream mode (i.e., when Alter3 fails to imitate with the automatic mimicry capacity and it generates a motion from its memory). We also found that people sometimes imitate Alter3, even when Alter3 was in the awake mode (i.e., when it generates a motion based on its autonomous mimicry capacity). It is interesting that people try to imitate the posture of Alter3 because it shows that imitation is an essential property of a living system. In other words, as people grow up, primitive imitation behavior does not disappear, but exists as a background process.

For example, close friends are known to synchronize the timing of their speech.

Starting from primitive imitation without any memories, Alter3 develops its memories via imitating human behavior and generates various behaviors based on memory selection and variation processes. While Alter3 interacts with a human and fails to imitate the human's behavior, humans tend to imitate Alter3 instead. This is quantified by the reversal of TE. We say that this spontaneous switching of roles between man and machine is a necessary condition of personogenesis.

DATA AVAILABILITY STATEMENT

The datasets generated for this study are available on request to the corresponding author.

AUTHOR CONTRIBUTIONS

AM and NM contributed equally to this paper, as first authors. All the authors designed the study and the model. AM and NM developed the system, performed the experiment, and analyzed the data. All authors contributed to the interpretation of the results. All authors drafted and revised the article. TI supervised the project.

REFERENCES

- Bongard, J., Zykov, V., and Lipson, H. (2006). Resilient machines through continuous self-modeling. *Science* 314, 1118–1121. doi: 10.1126/science.1133687
- Bradski, G. (2000). *The OpenCV Library*. Dr. Dobbs' Journal of Software Tools.
- Cao, Z., Simon, T., Wei, S., and Sheikh, Y. (2017). "Realtime multi-person 2D pose estimation using part affinity fields," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (Honolulu, HI), 1302–1310. doi: 10.1109/CVPR.2017.143
- Dayan, P., Hinton, G. E., Neal, R. M., and Zemel, R. S. (1995). The Helmholtz machine. *Neural Comput.* 7, 889–904. doi: 10.1162/neco.1995.7.5.889
- Farneback, G. (2003). "Two-frame motion estimation based on polynomial expansion," in *Image Analysis. SCIA 2003. Lecture Notes in Computer Science*, Vol. 2749, eds J. Bigun, and T. Gustavsson (Berlin; Heidelberg: Springer), 363–370. doi: 10.1007/3-540-45103-X_50
- Frith, C. D., Blakemore, S.-J., and Wolpert, D. M. (2000). Abnormalities in the awareness and control of action. *Philos. Trans. R. Soc. Lond. Ser. B Biol. Sci.* 355, 1771–1788. doi: 10.1098/rstb.2000.0734
- Ha, D., and Schmidhuber, J. (2018). "Recurrent world models facilitate policy evolution," in *Advances in Neural Information Processing Systems*, Vol. 31, eds S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Curran Associates, Inc.), 2450–2462.
- Helmholtz, H. (1867). "Handbuch der physiologischen Optik," in *Allgemeine Encyclopädie der Physik*, Vol. 9, ed G. Karsten (Leipzig: Voss).
- Hinton, G. E., and Sejnowski, T. J. (1983). "Optimal perceptual inference," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (Washington, DC).
- Hussein, A., Gaber, M. M., Elyan, E., and Jayne, C. (2017). Imitation learning: a survey of learning methods. *ACM Comput. Surv.* 50, 21:1–21:35. doi: 10.1145/3054912
- Iizuka, H., and Ikegami, T. (2004). Adaptability and diversity in simulated turn-taking behavior. *Artif. Life* 10, 361–378. doi: 10.1162/1064546041766442

FUNDING

This work was partially supported by the MEXT project Studying a Brain Model based on Self-Simulation and Homeostasis (19H04979) in Correspondence and Fusion of Artificial Intelligence and Brain Science as a Grant-in-Aid for Scientific Research on Innovative Areas. The authors declare that this study received funding from mixi, Inc. The funder was not involved in the study design, collection, analysis, interpretation of data, the writing of this article or the decision to submit it for publication.

ACKNOWLEDGMENTS

The authors would like to thank Dr. Kohei Ogawa, and Hiroshi Ishiguro from Osaka University for valuable discussions and also technical assistance with the engineering of Alter3. The authors appreciate mixi, Inc., for developing the virtual simulator of Alter3.

SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/frobt.2020.532375/full#supplementary-material>

- Ikegami, T. (2010). "Studying a self-sustainable system by making a mind time machine," in *S3'10: Workshop on Self-Sustaining Systems* (ACM), 1–8. doi: 10.1145/1942793.1942794
- Ikegami, T. (2013). A design for living technology: experiments with the mind time machine. *Artif. Life* 19, 387–400. doi: 10.1162/ARTL_a_00113
- Ikegami, T., and Iizuka, H. (2007). Turn-taking interaction as a cooperative and co-creative process. *Infant Behav. Dev.* 30, 278–288. doi: 10.1016/j.infbeh.2007.02.002
- Izhikevich, E. M. (2003). Simple model of spiking neurons. *IEEE Trans. Neural Netw.* 14, 1569–1572. doi: 10.1109/TNN.2003.820440
- Kingma, D., and Welling, M. (2014). Auto-encoding variational bayes. *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14–16, 2014, Conference Track Proceedings, (2014) arXiv [Preprint] arXiv:1312.6114*.
- Kwiatkowski, R., and Lipson, H. (2019). Task-agnostic self-modeling machines. *Sci. Robot.* 4:eaa9354. doi: 10.1126/scirobotics.aau9354
- Lizier, J. T. (2014). JIDT: an information-theoretic toolkit for studying the dynamics of complex systems. *Front. Robot. AI* 1:11. doi: 10.3389/frobt.2014.00011
- Masumori, A., Doi, I., Smith, J., Aoki, R., and Ikegami, T. (2020). "Evolving acoustic niche differentiation and soundscape complexity based on intraspecific sound communication," in *Artificial Life Conference Proceedings*, 465–472. doi: 10.1162/isal_a_00296
- McInnes, L., Healy, J., and Melville, J. (2018). UMAP: uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*. doi: 10.21105/joss.00861
- Meltzoff, A. N., and Moore, M. K. (1989). Imitation in newborn infants: exploring the range of gestures imitated and the underlying mechanisms. *Dev. Psychol.* 25, 954–962. doi: 10.1037/0012-1649.25.6.954
- Nadel, J., Revel, A., Andry, P., and Gaussier, P. (2004). Toward communication: First imitations in infants, low-functioning children with autism and robots. *Interact. Stud.* 5, 45–74. doi: 10.1075/is.5.1.04nad
- Nehaniv, C. L., and Dautenhahn, K. (2007). *Imitation and Social Learning in Robots, Humans and Animals: Behavioural, Social and*

- Communicative Dimensions*. Cambridge: Cambridge University Press. doi: 10.1017/CBO9780511489808
- Noguchi, W., Iizuka, H., and Yamamoto, M. (2019). Navigation behavior based on self-organized spatial representation in hierarchical recurrent neural network. *Adv. Robot.* 33, 539–549. doi: 10.1080/01691864.2019.1566088
- O'Keefe, J., and Nadel, L. (1978). *The Hippocampus as a Cognitive Map*. New York, NY: Clarendon Press; Oxford University Press Oxford.
- Piaget, J. (1966). *Play, Dreams and Imitation in Childhood*. New York, NY: W. W Norton & Co.
- Poincaré, H. (1905). *La valeur de la science*. Paris: Flammarion.
- Rössler, O. E. (1981). An artificial cognitive map system. *Biosystems* 13, 203–209. doi: 10.1016/0303-2647(81)90061-7
- Rössler, O. E. (1987). Chaos in coupled optimizers. *Ann. N. Y. Acad. Sci.* 504, 229–240. doi: 10.1111/j.1749-6632.1987.tb48735.x
- Rössler, O. E., Vial, L.-R., Kuske, F., Nitschke, A., Ikegami, T., and Ujica, A. (2019). Brain equation and personogenesis. *Clin. Pediatr.* 2:1011.
- Schaal, S. (1999). Is imitation learning the route to humanoid robots? *Trends Cogn. Sci.* 3, 233–242. doi: 10.1016/S1364-6613(99)01327-3
- Schreiber, T. (2000). Measuring information transfer. *Phys. Rev. Lett.* 85, 461–464. doi: 10.1103/PhysRevLett.85.461
- Suzuki, A., Josselyn, S. A., Frankland, P. W., Masushige, S., Silva, A. J., and Kida, S. (2004). Memory reconsolidation and extinction have distinct temporal and biochemical signatures. *J. Neurosci.* 24, 4787–4795. doi: 10.1523/JNEUROSCI.5491-03.2004
- Tani, J. (1996). Model-based learning for mobile robot navigation from the dynamical systems perspective. *IEEE Trans. Syst. Man Cybernet. B* 26, 421–436. doi: 10.1109/3477.499793
- The AuRoRA Project (1998). *The Aurora Project*. Available online at: <http://aurora.herts.ac.uk/> (accessed July 10, 2020).
- Trevarthen, C. (1977). "Descriptive analyses of infant communicative behavior," in *Studies in Mother-Infant Interaction*, ed H. R. Schaffer (London: Academic Press), 227–270.
- Wamsley, E. J., Tucker, M., Payne, J. D., Benavides, J. A., and Stickgold, R. (2010). Dreaming of a learning task is associated with enhanced sleep-dependent memory consolidation. *Curr. Biol.* 20, 850–855. doi: 10.1016/j.cub.2010.03.027

Conflict of Interest: AM, NM, and TI are employed by Alternative Machine Inc.

Copyright © 2021 Masumori, Maruyama and Ikegami. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

Advantages of publishing in Frontiers



OPEN ACCESS

Articles are free to read for greatest visibility and readership



FAST PUBLICATION

Around 90 days from submission to decision



HIGH QUALITY PEER-REVIEW

Rigorous, collaborative, and constructive peer-review



TRANSPARENT PEER-REVIEW

Editors and reviewers acknowledged by name on published articles

Frontiers

Avenue du Tribunal-Fédéral 34
1005 Lausanne | Switzerland

Visit us: www.frontiersin.org

Contact us: frontiersin.org/about/contact



REPRODUCIBILITY OF RESEARCH

Support open data and methods to enhance research reproducibility



DIGITAL PUBLISHING

Articles designed for optimal readership across devices



FOLLOW US

@frontiersin



IMPACT METRICS

Advanced article metrics track visibility across digital media



EXTENSIVE PROMOTION

Marketing and promotion of impactful research



LOOP RESEARCH NETWORK

Our network increases your article's readership