# An AI Theory of Mind Can Enhance Our Collective Intelligence

Michael S. Harré, Catherine Drysdale, Jaime Ruiz-Serra

July 9, 2025

## Abstract

*Collective intelligence* plays a central role in many fields, from economics and evolutionary theory to neural networks and eusocial insects, and is also core to work on *emergence* and *self-organisation* in complex-systems theory. However, in human collective intelligence there is still much to understand about how specific psychological processes at the individual level give rise to self-organised structures at the social level. Psychological factors have so far played a minor role in collective-intelligence studies because the principles are often general and applicable to agents without sophisticated psychologies. We emphasise, with examples from other complex adaptive systems, the broad applicability of collective-intelligence *principles*, while noting that *mechanisms* and *time scales* differ markedly between cases. We review evidence that flexible collective intelligence in human social settings is improved by a particular cognitive tool: our Theory of Mind. We then hypothesise that AIs equipped with a theory of mind will enhance collective intelligence in ways similar to human contributions. To make this case, we step back from the algorithmic basis of AI psychology and consider the large-scale impact AI can have as agential actors in a "social ecology" rather than as mere technological tools. We identify several key characteristics of psychologically mediated collective intelligence and show that the development of a Theory of Mind is crucial in distinguishing human social collective intelligence from more general forms. Finally, we illustrate how individuals, human or otherwise, integrate within a collective not by being genetically or algorithmically programmed, but by growing and adapting into the socio-cognitive niche they occupy. AI can likewise inhabit one or multiple such niches, facilitated by a Theory of Mind.

## 1 Introduction

*All intelligence is collective intelligence.* [1]

Collectives are capable of achieving things that individuals alone cannot. Notwithstanding the simplicity or complexity of the individuals, their aggregate behaviour can often be understood as a complex processing of information that individuals store, modify, and transfer between each other producing 'useful' collective behaviour at the scale of the whole collective. In most instances of *Collective Intelligence* (CI), where the agents might be ants in an ant colony, bees in a beehive, or neurons in a neural network, the individual is not aware of the drivers of their behaviour or the behaviour of other agents. For example, a single neuron is neither aware of its own internal processes nor that of a neuron it is connected to, nor is it aware of the end goal to which its activity contributes. Despite both this lack of awareness and the lack of a centralised controller, evolutionary and learning processes have produced an intricate, precise, and highly adaptive system that is capable of functional behaviour that would be impossible for any single neuron to achieve. In other instances of CI, such as teams of humans, or businesses interacting in economic markets, the agents themselves may be highly complex and exhibit varying degrees of purposefulness and awareness. Within this context, we draw attention to the role of psychological factors in improving the CI of human social collectives and quantifying the intelligence of social collectives, both natural and artificial. In this Introduction we review key aspects of collective intelligence and recent progress in placing AI into the collective intelligence framework.

### 1.1 Individual agents and their collective intelligence

In order to understand how collectives process information, we first consider the variety of ways in which agents interact. The topology of the network describing agent-to-agent interactions is well known to be important for the proper functioning of social groups [2, 3]. In particular it has been shown that mammalian social groups exhibit patterns of fractal-like topologies [4, 5] that are a result of a cognitive ability to form discrete social connections between conspicifics [6]. These links are often both spatially and temporally transient; people meet for a while, go their separate ways, and come back together later. Despite this transience, individual connections are often the basis of long term social relationships between specific individuals as in pair-bonding and friendships. Consequently an important distinction can be made regarding connections between agents in complex adaptive systems: they can be more fluid-like or more solid-like [7]. For example the links between neurons in the brain are relatively fixed in nature when compared to the brief communicative interactions between ants, either instantaneous interactions between individual ants or via transient pheromone trails that

coordinate the behaviour of large numbers of ants. Solé and colleagues [7, 8] identify a distinction between solid brains, in which interactions between agents fixed in place are highly persistent in time (e.g. neural networks, spin glasses) and liquid brains, in which interactions between highly mobile agents are much more short-lived (e.g. ants, immune cells). As Solé *et al.* note regarding liquid brains [7]: "Here there are no neural-like elements and yet in many ways these systems solve complex problems, exhibit learning and memory, and make decisions in response to environmental conditions."

All biological agents are composed of sub-units such as organs, cells, and molecular networks [9, 10, 11]. Cells in particular are the simplest living organisms with individual intelligence, or *competencies* [9, 12], within their native contexts. Here, we briefly focus on the archetypal single-cell intelligence, the neural cells. It is well understood that the central nervous system is a highly developed, adaptive, complex system that exhibits emergent computational characteristics [13], both in biological and artificial neural networks. Naturally the artificial models are simplifications but the extent to which they are simplifications is not so well understood. In a 2021 study, Beniaguev *et al.* [14] concluded that between five and eight layers of an artificial deep neural network are required to approximate the input–output mapping of a (single) cortical neuron and that the dendritic branches can be understood as spatiotemporal pattern detectors. This demonstrates that a single neural cell can be modelled as an artificial agent with highly complex computational capabilities situated within an adaptive, complex network of other highly complex agents, all signalling to one another. These results can be compared with earlier studies in which neurons were modelled as a Bayesian agent that is trying to infer the state of a hidden variable [15]. In each of these interpretations, a single cell can be seen as an agent with computational competencies situated within the context of a network that is slowly and adaptively changing around it.

We can also compare the competencies of neural cells in networks to the individual competencies of ants in an ant colony. In a recent study [16] it was shown that social structures of some ant colonies are conserved between species that are separated by more than 100 million years of evolution. In the five species studied by Kay *et al.* [16], they found two social clusters and similarities in the division of labour that are preserved between the species. In a different study, Richardson *et al.* [17] showed that individuals within an ant colony play an important *leadership* role and that the behaviour of these individuals significantly improved the collective performance of the ants. Ants are also capable of changing their social structure in the event of pathogenic infestation of their colony. In a 2021 article, Stockmaier and colleagues [18] review the research on *social distancing* and other social restructuring that occurs with conspecifics in order to reduce the impact of pathogens by changing their social cues, signals, and other behaviours for the collective benefit of the colony. These two very different systems, neural networks and ant colonies, are examples of complex collective intelligences where the individuals (neurons, ants) are complex in their own right, but they signal each other in order to restructure their relationships so as to adapt their collective competencies to external signals. The neural networks are prototypical *solid brains* and ant colonies are prototypical *liquid brains* and there has been recent progress in developing collective intelligence in the context of collective adaptation [19].

Human social interactions can also be viewed as a form of liquid intelligence. Migliano *et al.* [3] discuss the 'fluidity' of social relations in early human societies: "Quantification and mapping of hunter–gatherers' social networks has revealed details of a fluid and multilevel sociality, where friendship links connect unrelated mobile households into camps of temporary composition". They describe the key characteristics of early human society, such as egalitarianism, division of labour, cooperative living with unrelated individuals, multi-locality, fluid social structures, and high mobility between campsites, which might be thought of as a liquid brain composed of social interactions that both cluster and disperse in order to store, modify, and transfer information via social networks. The notion that human social interaction might be a form of computation is not new: Mirowski, Axtell and colleagues [20, 21, 22] have suggested that economic markets are a form of computation by which prices can be derived, and Harré recently hypothesised [23] that this could be measured using information theory as had been done earlier for financial markets [24, 25]. As Axtell *et al.* [21] wrote: "There is a close connection between agent computing in the positive social sciences and distributed computation in computer science, in which individual processors have heterogeneous information that they compute with and then communicate to other processors."

The emergence of computation in multi-agent systems is a well-studied area of complex adaptive systems [26, 27]. For example neuroscience has used information theory to describe the storage, transfer, and modification of bits of information in biological neural processes [28]. More broadly, Integrated Information Theory (IIT) [29, 30] has been put forward as a measure of the emergence of 'consciousness' in generic (non-biological, non-neural) systems. In this case, some forms of IIT explicitly use information theory [31, 32] to measure the amount of non-trivial computation a system is carrying out. More generally, there is a move towards understanding complex adaptive systems in computational terms [33, 34] by empirically measuring the inter-agent flow of information [35].

## 1.2   Recent developments using AI agents in complex adaptive systems

At the meso-scopic level, between artificial neurons and the collective behaviour of AI, is the emerging discpline of *machine behaviour* [36]. This field studies intelligent machines not as engineering artefacts but as a distinct class of actors with unique behavioural patterns and ecological dynamics [37]. Motivated by the increasingly pervasive integration of

algorithms into human society and the challenges in formalising and predicting their complex effects, this field examines machine behaviours empirically, similar to how ethology and behavioural ecology study animals by integrating intrinsic and environmental influences [38]. In this context, it bridges micro- and macro-level questions about algorithmic interactions and their societal implications, using methods like randomised experiments and observational studies from the behavioural sciences. By drawing parallels with biological frameworks, including Tinbergen's four dimensions of function, mechanism, development, and evolution [39], machine behaviour explores how machines interact with their environments and stakeholders, creating novel trajectories of evolution and influence that differ from organic systems. This interdisciplinary approach highlights the necessity of understanding AI's dual role in shaping and being shaped by human systems, suggesting careful examination of its impact to harness its benefits while mitigating risks in a socio-technical context.

We can illustrate the current state of specific AI frameworks in complex adaptive systems using Large Language Models (LLMs) as an example [40]. LLMs are powerful tools for enhancing CI, but they do not yet function as autonomous agents in their own right, but instead act as mechanisms that augment human capabilities in collaboration, creativity, and decision-making. Their potential applications are highly diverse, including increasing accessibility and inclusion in online collaborations, accelerating idea generation, mediating deliberative processes, and aggregating information across groups. However, their use also introduces risks, such as disincentivising contributions to collective knowledge, fostering illusions of consensus, reducing functional diversity, and enabling the seamless production of false or misleading information. Viewed with a wide lens, LLMs not only support CI but are themselves products of it, having been trained on collective human data and refined through collective feedback. This places the use of LLMs in the context of a tool that can improve pre-existing human CI, but the LLM itself does not have agential properties in the sense that its behaviour is not produced by selecting from multiple, internally developed goals, it has very little autonomy, and it is not typically connected to a real-time working environment.

In contrast to LLMs as non-agential AIs, new work is developing cognitively informed approaches to AI-human cooperation in which they have multiple self-regarding goals that account for the goals of others, allowing them to behave with a degree of cooperative agency. In recent work by Crandall and colleagues [41], it was shown that general human–machine cooperation is achievable through a sophisticated yet fundamentally simple set of algorithmic mechanisms, with results showing that an algorithm can cooperate with humans and other machines at levels comparable to human–human cooperation in repeated stochastic games (RSGs). Through a comparison of existing algorithms, the development of a novel learning algorithm that integrates high-performing strategies with human-conducive signaling mechanisms, the algorithm (S#) was shown to form and maintain cooperative relationships across diverse RSGs. Key to its success is the inclusion of 'cheap talk' [42] which significantly increases mutual cooperation and highlights the importance of signaling mechanisms that are aligned with human understanding. S# also employs a variety of expert strategies that can be adapted to different partners and game types. Its expert-selection mechanism, modelled on the psychologically based recognition-primed decision-making [43], allows S# to achieve a level of flexibility and generality unmatched by traditional expert-selection methods, illustrating its potential for fostering effective cooperation in complex, real-world interactions.

Similar results have been achieved in games with more realistic complexities, approaching human levels of subtlety in social reasoning within multi-agent, mixed human-AI, environments. Work carried out on an AI playing the game *Diplomacy* [44] has shown that AIs can successfully defeat expert human players in complex social games of cooperation and deception without being recognised as an AI. The AI, called Cicero, integrates dialogue generation, strategic reasoning, and message filtering into its decision-making architecture in order to pass messages to other players (all of them humans who were unaware Cicero is an AI) in order to strategise about future plans in the game. Its dialogue system is based on a pre-trained language model that is further developed on human game dialogue, grounded in both the dialogue history and the game state. Dialogue generation is controlled through 'intents'—planned actions involving Cicero and the person it is messaging—enhancing contextual relevance and strategic alignment while offloading the responsibility of learning game legality and strategy to other modules. Cicero's strategic reasoning module employs a planning algorithm that predicts opponents' policies based on game state and dialogue, selecting optimal actions guided by reinforcement learning with constraints to align with human-like behaviour. The approach used by Cicero has been described as an elementary form of modelling the cognitive state of other players, i.e. having a Theory of Mind. Messages are dynamically generated and filtered to ensure they are coherent, consistent with intents, and strategically sound, allowing Cicero to integrate negotiation and strategic planning in a manner indistinguishable from human players. As one of the researchers describes it [45]: "I think that the substance of what we're doing is understanding the beliefs, goals and intentions of the other players.". Both Cicero and S# illustrate that language combined with psychological aspects of reasoning and a strategic awareness of the environment are able to build coordination and cooperation with humans well enough for the AI to pass as a human. In order to move towards AI based more closely on human psychology [46] to enhance our collective intelligence [47] requires a greater awareness of the ecology and sociology of AI [48], and the development of interdisciplinary approaches to understanding the macroscopic consequences of these technologies.

In this article we will use information theory to quantify how much computation in a CI is 'emergent' and how much is simply independent information processing by single agents. In general, we wish to capture the notion of the whole (computational process) being greater than the sum of the (independent) parts. We translate this to the simple notion that to the extent to which this inequality holds: `Whole - ∑(Parts)` $> 0$ is the extent to which we will say a system

exhibits non-trivial CI, noting that there are multiple possible implementations of this approach [49]. The `Parts` is how much computation a single agent is carrying out from one time step to the next such that the sum is the total of all agents' independent computations. The `Whole` is the totality of computation in the system, it includes all single agent computations, pairwise computations, and higher order interactions between agents. Our measure will not be unique in any of its specifics, but it serves to quantify the CI of a system for comparative analysis. This approach also has much in common with that of Moore *et al.* [27] in which information theory is used to measure the collective intelligence in biological systems.

Not only is there diversity in the types of systems that can show positive measures of CI, but the ways in which agents manipulate a system's computations is diverse as well. Take for example Watson and Levin's discussion of a scientist manipulating the intercellular signalling in order to change their collective outcome [50]:

> This framework [of collective cellular intelligence] makes a strong prediction: if intercellular signalling (not genes) is the cognitive medium of a morphogenetic individual, it should be possible to exploit the tools of behavioural and neuro-science and learn to read, interpret and re-write its information content in a way that allows predictive control over its behaviour (in this case, growth and form) without genetic changes.

A counter question is: How can single agents, such as human leaders, have predictive control over a social group? Just as a scientist *external* to a cell collective can manipulate inter-cellular signalling to control the outcomes of the cell collective, a leader *internal* to a human collective can manipulate inter-personal behaviours to control the outcomes of the human collective. In both cases, an agent with a goal-directed psychology is acting on inter-agent relationships, i.e. inter-cellular or inter-personal, to control outcomes at the next level higher, i.e. organism-scale or societal-scale.

## 1.3   The approach taken in this article

In this article we will reframe the question Watson and Levin asked in the following way: *What is there in human psychology that allows us to learn to read, interpret, and re-write our interpersonal information content in a way that allows predictive control over our collective behaviour?* and consider the answer in the context of pervasive, agential AI. We review some of the extensive literature showing that our Theory of Mind (ToM) is a suite of cognitive skills that allows individuals to have goal directed control over collective outcomes. Originally ToM was used to describe our ability to infer the unobserved mental states of other people [51] such as desires and beliefs, an ability humans are particularly good at and other animals much less so [52, 53]. But recently it has been shown that ToM is predictive of group performance as well [54, 55], empirically demonstrating the role of ToM in going beyond representations of the internal states of others to using that knowledge in a social setting to improve the collective outcomes for the group. In order to model ToM in a tractable fashion, we will focus on the narrower *game theory of mind* [56], and the *Beliefs, Preferences, and Constraints* (BPC) interpretation of game-theoretic decisions put forward by Gintis [57]. In this approach, what agents understand of other agents' hidden states are the BPC that structure their observable behaviours.

We will consider this question in the framework of agent interactions that extend agent utilities in a simple but novel way. We quantify our results using information theory to show the impact that a well-developed ToM has to direct agents' behaviours in order to increase our CI. The models are simple but they illustrate the central notion that understanding the "beliefs, preferences, and constraints" [58, 59] of others can be used to improve the CI of a complex social system. The wider purpose of this work is to place recent developments, such as human-AI co-evolution [60], AI-enhanced collective intelligence [47], the collective sociology of AI and humans [48], the connection between ToM and socio-cultural niche evolution [61], and the design of intelligent cyber-physical ecosystems [62], into the context of measuring hybrid CI, CI facilitated human sociology via ToM, and how concepts from ecology, psychology, and economics can help us build a better understanding of the joint future evolution of humans and AI.

In Section 2, we describe the liquid–solid spectrum of interacting agents, review existing models of ToM, and provide perspective on the interplay between social network structures and ToM. We specifically draw attention to three key ideas that provide empirical support for either ToM or CI at the scale of human collective behaviour. First, the liquid–solid brain hypothesis establishes a cognitive spectrum from rigid and persistent to fluid and dynamical interactions between agents that exhibit collective intelligence. This approach is independent of the spatial scale of the collective, for example it can be used to study neurons, insect colonies, or human social interactions. Second, we consider the evidence for ToM, its structure at the individual agent level as a method for agents to model the unobserved control parameters that influence other agents' behaviour. We then show how our use of language, a fluid communication channel that temporarily couples agents together in order to share information and modify behaviour, is tightly connected with our development and use of ToM. Finally, we consider the role of network topology in CI and in particular how recent work has measured the CI of social groups and related it to the individual's capacity for ToM. In Section 3, we provide illustrative examples supporting different aspects of our argument, introducing our measure of computation and applying it to a simple empirical example. Here we note that the effectiveness of a collective depends on the macroscopic structure of interactions between agents such that any analysis will need to account for more than simple dyadic interactions and so we introduce these higher order

interactions. In computing the CI between agents, a computer interacting with a monkey, we illustrate a game theory environment with an AI-biological hybrid system in which the CI value can be measured using observable behaviours. In Section 4 we review the psychology of social fluidity and the variety of social outcomes that this fluidity makes possible. We also use a simple multi-agent system to describe how a ToM can be used to improve the computational processes, i.e. the CI, of interacting agents. It is here that we connect the hidden variables (i.e. meta-parameters) argument of ToM discussed earlier with the parameters of game theory describing the beliefs, preferences and constraints that incentivise decisions [63, 64]. We also provide a formal interpretation of the higher order interactions introduced earlier in terms of game theory. Finally, in Section 6, we discuss the broader implications of this approach.

# 2 Cognitive morphospaces: The network topology of agents' interactions

The emergence of cognitive networks marked a pivotal moment in our evolutionary history [65]. Earlier, microorganisms had developed collective structures capable of responding to the physical environment, particularly conditions threatening individual cells [66] and survival became dependent on information exchanges within these groups. Habituation and the ability to minimise the energy response to danger stimuli is considered one of the simplest forms of learning and has been extensively studied in simple collective intelligences such as slime moulds [67, 68]. Similarly bacteria use quorum sensing to coordinate their behaviour based on the density of the population of their community to coordinate responses, leading to a change in gene expression and function regulation e.g. bioluminescence, release of toxins, and biofilm formation [69]. Information processing and problem solving capabilities developed using a variety of interaction types and network topologies long before the appearance of central nervous systems with fixed neuronal structures [66]. But this leads to an interesting question regarding the *typologies* of agent-to-agent interactions and the intelligence these structures might enable, a question that can be approached by looking at the *morphospace of collective intelligence.*

A morphospace is a theoretical framework used to simplify and organize the complex shapes and forms of organisms, typically focusing on external anatomical features, into a more manageable space representing their potential variations. For instance, in [70], a three-dimensional morphospace was constructed for organisms with shells, where the diversity of shell shapes is described by three key parameters: a deviation angle, a translation factor, and a growth factor. This reduces a high-dimensional structural space to a lower-dimensional one, where a small number of parameterised properties captures key variations between forms. Morphospaces have found uses in many different fields of research, including the body shapes of fish [71], network topologies [70], and the structural forms of language [72].

These structures, which represent recurring patterns of trait variation, are of great interest to evolutionary biologists because they may indicate shared evolutionary processes and their constraints. They are also of great interest to researchers investigating the underlying structures of collective intelligence [73], as intelligence in its different forms may also be subject to shared evolutionary mechanisms and constraints. With this in mind, work has been done in studying a variety of morphospaces related to collective intelligence and in the next section we review some examples.

## 2.1 Network topology and the "Solid Brain, Liquid Brain" framework

In the most general of terms, a cognitive network has multiple information processing components that exchange information with each other and interact with their environmental context. Solé *et al.* [7] identified two key dimensions for categorising different types of cognitive networks: the system's physical state—either more *liquid* or more *solid* in nature, differentiated by how freely individual agents (components) can move in space—and the presence or absence of neurons. The collective dynamics of a large population of agents is influenced by the individuals' mobility which dictates how they respond to signals both internal and external to the collective. To help conceptualise this diversity, Solé *et al.* [7] developed a morphospace and taxonomy in order to compare and contrast the physical states of different types of CI. In this way they were able to consider the physical properties that form constraints on the computations achievable by a system. This then poses an interesting question: Is the entire space of possible *intelligences* being exploited by either synthetic biology or abiotic computation?

Liquid brains exhibit cognitive behaviours without neurons. For example, models by Watson and colleagues [74, 75] illustrated how systems of self-interested agents, driven by a simple mechanism of strengthening beneficial connections, can lead to robust group-level adaptation and problem-solving. This self-organisation, akin to Hebb's rule in neural networks, enables the system to recall and consequently leverage past configurations that were successful. This also allows the system to generalise from experience and then predict beneficial states it has not encountered before, highlighting how decentralised actions can produce a form of CI that guides the system towards greater global utility. The agents in these models are very simple, but this need not be the case—each agent within a collective may itself have a solid brain as in human social networks where fluid social interactions allow each solid brain to connect and communicate with other solid brains for the benefit of the collective [76]. This is extended to a hybrid model by Kao *et al.* [77], wherein the modular organisation of mobile animal populations (as an example of a liquid brain) suggests that certain communication

pathways exhibit localised and persistent characteristics, akin to those in solid brains. These pathways enhance collective decision-making in complex environments and in turn raise important questions regarding the relative strengths of liquid and solid brains in different contexts. In particular, the conditions under which liquid brains outperform solid brains, especially in terms of adaptability and scalability, remain an open area of investigation. Which computational problems are more effectively solved by liquid brains vis-à-vis solid brains remains to be better understood.

Whole classes of models that describe which biological or artificial structures are capable of some form of computation, either in potentia or in practice, can often be usefully represented using morphospaces [78]. Computational morphospaces [7] have proven effective in studying key properties of complex adaptive systems, for example the statistical mechanics of information processing and structural variations [79, 80]. This sheds light on how energy constraints influence the evolution and adaptability of neural networks across a variety of biological systems. Arsiwalla *et al.* [79] have examined how liquid brains and solid brains fit within such a framework, by comparing the flexibility and adaptability of different neural architectures [79]. The dimensions of their framework are three different types of complexity that a system may display: autonomic, computational, and social complexity. We suggest that other axes for consideration are a system's solid–liquid dimension as Solé *et al.* [7] and Ollé-Vila *et al.* [81] have done, as well as the system's degree of ToM (see Section 2.2) and the system's information processing capacity (see Section 3).

We posit that the collective intelligence that emerges from liquid brains (human social networks) is enhanced by our individual capacity for a ToM, where individuals are aware of the goals of others as well as that of the collective, which can then be achieved by adaptation at the local level. Rather than collective intelligence arising as an epiphenomenon or byproduct of agents interacting, our ToM allows agents to causally affect the outcome of the system they are a part of.

## 2.2 Models of another agent's internal states

Frith and Frith defined Theory of Mind (ToM) as how we explain other people's behaviour on the basis of their internal cognitive states, i.e. their knowledge, beliefs, and desires [51]. There is now a vast literature on this topic in psychology, sociology, and more recently artificial intelligence, but for the purposes of this article we restrict ToM to apply to the subset of the beliefs, preferences, and constraints of other agents in the sense of incentivised decisions. This borrows from the BPC model [57] put forward as an approach to understanding the socio-cognitive aspects of human decision-making [59].

One way to interpret the BPC model is that it imposes structural constraints on the process by which decisions are made, and then optimal decisions are discovered within these constraints by parametric variation. Recent work by Peterson *et al.* [82] compared more than 20 structurally constrained models of individual decision-making using human data. That study was extended to human data during strategic interactions by Harré and El-Tarifi [63] in order to test agents' constrained representations of other agents. This extends Yoshida *et al.*'s [56] notion of a *Game Theory of Mind* to a larger variety of models in which an agent's strategic reasoning about other agents modulates their behaviour.

In the field of artificial intelligence, Jara-Ettinger [83] proposed the use of Inverse Reinforcement Learning (IRL) as a model of agential ToM, whereby agents modelling other agents' mental states is equivalent to inferring an unobserved world model the other agent uses in their decision-making, as well as their reward function. Jara-Ettinger discusses some key limitations of IRL, such as the difficulty in recovering an agent's beliefs and desires even while assuming that all agents are identical in their choice-making. IRL has been implemented in many different algorithmic forms, and their applicability as a basis for ToM was recently reviewed by Ruiz-Serra and Harré [84] and recent work by others in developing ToM for artificial intelligence in collaboration with people [85, 86, 87].

What is missing in the Jara-Ettinger perspective is how different internal models, specifically different drivers of behaviour predicated on the BPC of others, influence how agents interact with one another in a social network. Critically, people use and improve upon these information carrying interactions in their social networks, as discussed next. This requires agents to have more than a representation of the world model used by other agents, it needs to be a *social* world model of how agents are influenced by their interactions with other agents. Beyond inferring internal states, Shteynberg *et al.* [88] distinguish between awareness of the self, awareness of the self in relation to others, and an awareness of the collective, which is more than just the sum of the self and others, a type of collective awareness they have called *Theory of Collective Mind*. In a similar vein, Shum *et al.* [89] proposed a generative model for understanding multi-agent interactions called Composable Team Hierarchies. This approach used stochastic games in conjunction with multi-agent reinforcement learning in order to infer relationships between agents and predict future behaviours.

In cognitive science, recent progress has been made in identifying levels of ToM ability and placing them in a hierarchical structure (see also Yoshida *et al.*'s Game Theory of Mind [56] and a recent review by Harré [64]). A cognitive agent with *zeroth order* ToM attributes no cognitive ability to other agents, whereas *first order* ToM attributes some cognitive abilities to others, and so on. Here we summarise the orders as described by Lombard and Gärdenfors [90] and use *A* and *B* to identify two agents that may or may not have any cognitive ability:

- Zero order ToM: Both *A* and *B*'s behaviour is governed by instinct, reflexes, and conditioning and so direct perception of the agents' interactions with their environment and each other is all that is needed to understand their behaviour.

- First order ToM: *A* and *B* can be attributed with emotions, attention, desires, intentions, or beliefs, but neither agent attributes these properties to any other agent, including themselves.

- Second order ToM: *A* attributes to *B* internal cognitive states, and that *A* uses this knowledge to understand *B*'s behaviour. This is the lowest level at which *A* attributes hidden (cognitive) variables to *B* in order to explain the causes of *B*'s actions, i.e. it abstracts causation away from direct perception of the causes of behaviour.

- Third order ToM: *A* attributes to *B* an understating of *A*'s internal states. To borrow an example from game theory (stag hunt) and early human society (hunter gatherers), when a hunting party stalks an animal everyone shares a common goal whereby each person *A* knows that the others are aware of *A*'s goal, such that this cognitive state will causally inform *A*'s actions. Lombard and Gärdenfors note that it has not been conclusively demonstrated in nonhuman primates and indicate there are alternative views [91, 92].

- Fourth order and higher ToM: *A* has an awareness of at least two mental states, their own and that of *B*. For example Happé reviewed the evidence and suggested that reflecting on one's own cognitive state relies on the same neuro-psychological functions as those we use to attribute thoughts to others [93].

Most interactions between agents across all species will be of the zeroth or first order ToM, where neither agent has higher order cognitive states, has no sense of being aware of other agents nor any self-awareness. There is no sharp line that unambiguously distinguishes between stimulus-response, conscious, and self-conscious agents as it is an open area of research and is a multi-dimensional phenomenon [94], and so the conscious status of agents is likely more incremental than the discrete levels of this scale would indicate, but is a useful framing device.

This leads to another core finding related to ToM: the interplay between language and ToM in humans. Chomsky has noted both how well developed this ability is in humans and how communication sits in relation to our mental states [95, p. 10]:

> Communication is not a matter of producing some mind-external entity that the hearer picks out of the world, the way a natural scientist could. Rather, communication is a more-or-less affair, in which the speaker produces external events and hearers seek to match them as best they can to their own internal resources. Words and concepts appear to be similar in this regard, even the simplest of them. Communication relies on largely shared cognoscitive powers, and succeeds insofar as similar mental constructs, background, concerns, and presuppositions allow for similar perspectives to be reached. If that is true—and the evidence seems overwhelming—then natural language diverges sharply in these elementary respects from animal communication.

That is to say, effective communication between people requires an encoding-transmission step and a reception-decoding step, both premised on shared cognitive representations in order to be understood, which we contend is made possible via ToM, i.e. a (shared) representation of each other's hidden cognitive variables.

Recent work has informed Chomsky's view here, linking language and ToM with our ability to carry out causal reasoning, particularly in social contexts. In Lombard and Gärdenfors they proposed three key hypotheses relating ToM to cognitive structures [90] (quoted, emphasis added):

- Theory of mind is an integral element of *causal cognition*;

- Generally speaking, the more advanced causal cognition is, the more it is dependent on theory of mind; and

- The evolution of causal cognition depends more and more on *mental representations of hidden variables*. [...] [C]ausal cognition allows us to reason from a network of hidden variables ...

Language use is related to this causal reasoning about hidden variables via the representational view of language due to their intersection with ToM [96]. In the representational view of language people use specific grammatical structures to represent complex events and then to reason from them [97]. These structures serve as a cognitive tool, particularly in representing others' mental states, enabling the expression of false beliefs, lies, or mistakes. It has also been shown that these are strong predictors of children's false belief understanding, a canonical test of ToM, and targeted syntactical training can improve false belief reasoning [98]. Notably, in developmental learning, language has a stronger influence on ToM abilities than vice versa [99]. From this we identify causal reasoning about hidden cognitive states with causal reasoning about the behaviour of others via a shared representation of our social environment, mediated by specific syntactical structures.

How is language 'causal' though? That is, how do cognitive representations causally manifest themselves in collective outcomes via language? To some extent this might be a natural assumption to make, but its significance has been shown experimentally in that the words we choose, that grammatically represent our internal models of the world, have consequences in policymaking. In work by Thibodeau and Boroditsky [100] they have shown that two different metaphors used to describe the same numerical data regarding crime in a city: as either an infectious disease or a rampaging beast, causally influenced the likelihood of policies people chose: either preventive or punitive measures. This, they have

argued [101], is the conceptual scaffolding through which we reason and that these linguistic metaphors guide our thoughts and behaviours. To date there has been little work on the levels of ToM for AI or the causal nature of language models, and this gap in the literature will need to be filled if we wish to understand how AIs can play an effective role in human collective outcomes. But if language is causal, the degree to which our language impacts outcomes is also contingent on the structural properties of the networks over which this information spreads [102]. We discuss these aspects next in the context of shared cognitive representations.

## 2.3 Configurations of our social networks inform individual reasoning

Frith and Frith have previously considered the benefits that accrue to humans via our ToM [51], but what is functionally happening when we use our ToM in social groups? It has been shown that in early hunter-gatherer societies that some emergent phenomena at the social level, e.g. Dunbar's Number [103, 104], are a consequence of the layered, fractal topology of social networks [105, 4, 5], and that these are in turn the product of very specific, discrete, cognitive constraints at the individual level that shaped the structures of early human societies [6]. At the individual level, a recent review by Momennejad [2] collected the evidence for different social network topologies and how they integrate interpersonal knowledge differently, showing that topology allows social networks to serve a rich variety of collective goals. Momennejad also reviews the neuro-imaging evidence showing that humans neurologically encode these topologies and these encodings are shared across the members of a social group. This was also demonstrated in the work of Lau *et al.* [106] showing that people are able to integrate information about how agents relate to one another in addition to how they relate to oneself in order to infer social group structures. Lastly, there is evidence for improved collective intelligence when individuals with higher competencies in ToM are present in the group, as shown in the study by Woolley *et al.* [54]. It was found that, just as there is for an individual person a measure of general cognitive ability, usually denoted $g$, there is an equivalent measure of collective intelligence, denoted $c$, for a group of people. They noted a key explanatory factor of a group's task performance, as measured by $c$, was the proportion of group members who ranked highly on the *Reading the Mind in the Eyes* cognitive test introduced by Baron-Cohen and colleagues [107, 108], a test used to measure an individual's capacity for ToM.

A key takeaway from the work of Woolley and colleagues [54] is that the $c$ factor is not strongly correlated with either the average or maximum intelligence of the individuals in a group. However, it does correlate well with the average *social sensitivity* of its members, the evenness of the distribution of contributions to group discussions, as well as the proportion of people in the group who rate highly on a ToM test. So there is considerable evidence for the role our ToM plays in group performance, how this shapes interpersonal interactions between people, and, as a consequence, the emergent topological properties of our social groups. This provides support for the argument that ToM is (one of) the individual, bottom-up mechanism(s) through which agents form higher-order social structures with measurable collective intelligence. This may be how we "learn to read, interpret and re-write our interpersonal information content" as we asked following the Watson-Levin quote in Section 1. Next, we illustrate these ideas in a dyadic and then a triadic example of agents interacting with each other and exhibiting non-trivial measure of CI.

# 3 Topological and cognitive structures in CI: illustrative examples

## 3.1 A dyadic example of individual learning at short time scales

To illustrate how we will use information theory as a proxy for Woolley *et al.*'s $c$ intelligence, we use a dataset that was previously studied [109] to measure the information flow in an iterated economic game experiment between monkeys and computers, based on data from an earlier study by Lee *et al.* [110]. In that study, the experiment had a monkey playing the matching pennies game against a computer algorithm for a reward, the (Nash) optimal reward for the monkey was received if it plays 50:50 across its two choices. A simplified description of the three algorithms used by the computer follows, see Lee *et al.* [110] for the exact descriptions:

- Algorithm 0: Play uniformly and independently of the monkey's choices,

- Algorithm 1: The computer stores the history of the choices made by the monkey, then to predict what the monkey would do in each trial the computer calculates the conditional probability of the monkey's choice given the monkey's choices in the preceding 4 trials, a Null hypothesis was used to test if the monkey played 50:50 and if it was not rejected the computer plays 50:50, if it was rejected the computer plays probability $1 - p$ against the monkey's probability of playing $p$

- Algorithm 2: Uses the same algorithm as 1) but includes both choices and rewards in the monkey's history, and then both algorithm 1) and 2) were tested against the Null of 50:50 and if the Null was not rejected the computer plays 50:50, otherwise the best reply was played based on the estimated bias of the monkey.

The usefulness of this example is fourfold. First, it is simple enough that computations can be tested and evaluated in an illustrative way so that the central information theory measures can be applied. Second, it uses game theory as the foundational mechanism for the interactions that generate the information flow between the agents. Third, it is complex enough to illustrate key elements of CI with only two agents. Finally, it illustrates the general principles that can be applied in evolution, psychology, AI collectives, and human-AI hybrid settings. Our definition of CI is based loosely on the information theory form of Integrated Information Theory (IIT) [29], where we are neutral to the interpretation of IIT as a measure of consciousness, and will not be carrying out any optimisation over binary partitions of state variables, so in that sense this is not the same as IIT, although there are some similarities. Our definition for a system with $n$ agents is (see Section 2.2.3 of [32] and for a general introduction see Battencourt [111]):

$$\phi(X;\tau) \triangleq I(X_t; X_{t-\tau}) - \sum_{i=1}^{n} I(X_t^i; X_{t-\tau}^i). \tag{1}$$

For example we might have $X_t = \{X_t^1, X_t^2\}$ as the joint stochastic variable of agents 1 and 2 such as the monkey and the computer. The function $I(X_t^i; X_{t-\tau}^i)$ is the *time-delayed mutual information* (TDMI) with delay $\tau$ for any times series of a stochastic (possibly joint) variable $X_t \in \{X_1, X_2, \ldots, X_T\}$:

$$I(X_t; X_{t-\tau}) = \sum_{i=1}^{n} p(X_t, X_{t-\tau}) \log \left[ \frac{p(X_t, X_{t-\tau})}{p(X_t)p(X_{t-\tau})} \right]. \tag{2}$$

Equation 2 encodes the number of bits that variable $X_t$ is able to predict about its own future state based on its ($\tau$-lagged) past states. For example, in a two-agent system, if $X_t = \{X_t^1, X_t^2\}$ then $I(X_t; X_{t-\tau})$ is the amount of predictive information data $\tau$ steps in the past is encoded in the current state of the entire joint state of the system. On the other hand, if $X_t = X_t^1$ then $I(X_t; X_{t-\tau})$ encodes how much information the past of one agent encodes about its current behaviour. The difference between the whole system TDMI and the sum of the individual's TDMI is the extent to which information is being exchanged between the agents, i.e. $\phi(X;\tau)$ in Equation 1 is the *excess TDMI*.

Table 1 shows the results of these computations. We note the fact that, as the computer algorithm becomes more sophisticated, from algorithm 0 to 1 to 2, $\phi(X;\tau)$ decreases as the monkey plays closer and closer to the Nash strategy, but then as the sophistication increases, the monkey looks further into the past data to extract information. This is most notable in the difference for algorithm 2: at $\tau = 3$ (0.0148 bits) is nearly twice that of $\tau = 1$ (0.008 bits). This is because, in this specific case, the monkey decouples itself from the computer in order to earn its highest reward: as the computer strategy becomes more sophisticated in measuring the monkeys use of the computer's choices, the less coupled to the computer and its own past the monkey needs to be.

| Time delay: $\tau$ | Algo | Joint TDMI | Monkey TDMI | Computer TDMI | Excess TDMI: $\phi(X;\tau)$ |
|---|---|---|---|---|---|
| 1 | 0 | 0.0999 | 0 | 0 | 0.0999 |
| 2 | 0 | 0.1197 | 0.0075 | 0.0002 | 0.1120 |
| 3 | 0 | 0.1237 | 0.0095 | 0.0095 | 0.1047 |
| 1 | 1 | 0.0693 | 0.0011 | 0.0001 | 0.0681 |
| 2 | 1 | 0.0748 | 0.0025 | 0.0001 | 0.0722 |
| 3 | 1 | 0.0766 | 0.0029 | 0.0005 | 0.0732 |
| 1 | 2 | 0.0105 | 0.0005 | 0.0002 | 0.0080 |
| 2 | 2 | 0.0155 | 0.0019 | 0.0023 | 0.0113 |
| 3 | 2 | 0.0216 | 0.0030 | 0.0038 | 0.0148 |

Table 1: The simplest example of two agents interacting with one another via game theory in which $\phi(X;\tau)$ is non-zero. In general we cannot know the details of the flow of information, only that there is a net positive flow across all agents. All values are in bits and computations carried out using JIDT [112].

## 3.2   A triadic example of evolutionary learning at long time scales

Before we introduce ToM for social interactions, we consider a second example where evolution has found an agent-to-agent interaction that is similar to the worked example used next in Section 4. Our evolutionary example is a three-agent system: the larval stage of the fly *Liriomyza huidobrensis*, the pea plant family *Fabaceae* that fly larvae predate on, and the parasitic wasp *Opius dissitus* that predates on *L. huidobrensis*. These species interact in the following way [113]: The larvae of *L. huidobrensis* infest a pea plant, the pea plant gives off volatiles, called *infochemicals*, that attract the wasps to the plant, which in turn feed on the larvae, and this larva–wasp conflict indirectly benefits the pea plant.

In this triadic relationship, the pea plant does not directly respond to the threat from the larvae—for example it has not evolved a chemical agent that repels the larvae-laying flies. Instead it signals a third party, the wasp, to bring the wasp into contact with the larvae, and the wasp then eats the larvae. The wasps and the flies are in a dyadic evolutionary competition that could be modelled using two-agent evolutionary game theory, but the pea plant, having facilitated an instance of this conflict, benefits indirectly from the wasps success, in a sense plants can employ other species as a kind of 'body guard' [114]. We note that, like the example in Section 4, the plant only signals the wasps to a dyadic interaction when the plant detects the larvae, so that the plant only signals wasps when the plant perceives an intermittent information carrying cue from its environment that wasps are needed.

This type of *second order* interaction occurs in other ecological examples as well [115] where Trait-Mediated Indirect Interactions (TMIIs) induce hyper-graphs of interactions between species. More complicated interactions have also been observed in which a plant that is being attacked emits infochemicals that lead to *unaffected* plants emitting volatiles to attract predators [114], to reduce the risk of the unaffected plant being attacked while also aiding the plant being attacked. Kobayashi and Yamamura [114] specifically call this evolutionary development a form of *altruism* but of course there is no cognitive aspect to this altruism. These examples illustrate that evolution produces forms of sophisticated, inter-species, mixed competitive-altruistic interaction networks but without any need for a ToM, individual awareness, or strategic understanding of the interactions, and so individuals are not knowingly strategic or altruistic as we might interpret a person to be. But these evolved strategies are limited, by definition, to be fixed within the lifetime of a single agent and they can only adapt on evolutionary time scales. This is in contrast to a ToM which allows us to adapt to multiple strategic and social contexts that may require novel solutions within the lifespan of a single person.

n the following section we illustrate how the strategic modification of an interaction network leads to improved CI, contingent on hidden cognitive variables. This provides a formal connection between hidden variables and the causal consequences of network plasticity at the level of the collective. As this is a 'causal model' the agent is using, it provides a useful toy model of how an AI may use cognitive tools similar to those of a human to connect the causal use of language with the consequences of network adaptation for improved collective outcomes.

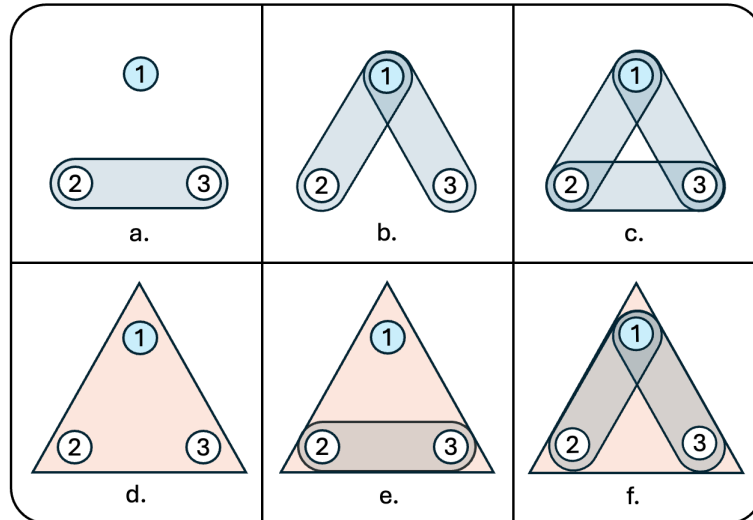# 4 Social network modification through ToM: a minimal model



Figure 1: Graphs and hyper-graphs of three agents: (a) A disconnected graph containing two dyadically connected agents and one isolated agent. (b) A connected dyadic graph. (c) A dyadic *complete graph*. (d) A hyper-graph in which all three agents are connected by a single hyper-link. (e) A combination of a hyper-graph connecting all agents and a disconnected dyadic graph. (f) A connected dyadic graph and a single hyper-link.

In this model we introduce a simple example of interacting agents that form a hypergraph [115]. As in the previous examples, this model is general enough to allow the agents to have any form of biological or artificial psychology, to have zeroth order ToM or fourth order ToM, and it applies to artificial collectives just as readily as it can to biological and ecological collectives.

Game theory provides a formal approach to the analysis of incentivised social interactions in which agents attempt to optimise the outcome (utility) of their joint actions. The expected utilities can be rewritten as polynomials in the agents' decision variables, usually interpreted as probabilistic weights, $x_i \in \mathbf{x}$ for which the co-factors $\mathbf{a}$ are derived from a game payoff matrix (as described in Appendix A). These can be generalised to higher order polynomials representing higher

order interactions between the agents: in these expanded utilities the quadratic terms represent the dyadic interactions between agents, the cubic terms represent the three-way interactions, and so on:

$$U_i(\mathbf{x}; \mathbf{a}) = a_i^0 + \sum_j a_i^j x_j + \sum_{j,k} a_i^{jk} x_j x_k + \sum_{j,k,l} a_i^{jkl} x_j x_k x_l + \text{h.o.t.} \tag{3}$$

We note that in general $\frac{\partial U_j(\mathbf{x};\mathbf{a})}{\partial x_i}$ describes the impact that $i$'s choice has on agent $j$'s utility and that we need not assume that there exist either symmetrical or even reciprocal impacts between agents' utilities and their behavioural choices. In Figure 1, we illustrate how three agents are connected via dyadic relationships and higher order (cubic) hyper-graph interactions. Taking the utility for Agent 1 in Figure 1 as an example, the most general description of the interactions in the dyadic utilities for the top row (a)–(c) are given by:

$$(a): \ U_1(\mathbf{x}; \mathbf{a}) = a_1^0 + a_1^1 x_1 \qquad\qquad (b)\,\text{and}\,(c): \ U_1(\mathbf{x}; \mathbf{a}) = a_1^0 + \sum_j a_1^j x_j + \sum_{j,k} a_1^{jk} x_j x_k \tag{4}$$

In (a), Agent 1 is the only agent that can influence their utility, but for (b) and (c), they are also influenced by other agents, linearly and possibly quadratically, depending on the payoff structure. The co-factors $\mathbf{a}$ of these three utilities are derivable directly from the utility matrices of dyadic agent-to-agent interactions in conventional game theory [116, 117] (see Appendix A). For utilities (d)–(f) in Figure 1:

$$(d): \ U_1(\mathbf{x}; \mathbf{a}) = \sum_{j,k,l} a_1^{jkl} x_j x_k x_l \qquad\qquad (e): \ U_1(\mathbf{x}; \mathbf{a}) = a_1^0 + a_1^1 x_1 + \sum_{j,k,l} a_1^{jkl} x_j x_k x_l \tag{5}$$

$$(f): \ U_1(\mathbf{x}; \mathbf{a}) = a_1^0 + a_1^1 x_1 + \sum_{j,k} a_1^{jk} x_j x_k + \sum_{j,k,l} a_1^{jkl} x_j x_k x_l \tag{6}$$

With these descriptions of higher order interactions we can now provide an illustrative example of how a ToM may be used to reconfigure a network to enhance the computational processing of information to improve collective performance. In the example of the parasitic wasp and the fly larvae eating the leaves of a plant, the plant benefits indirectly from the direct interaction between two other agents. In this case the plant encourages the wasp to come into contact with the larvae in a fight for survival that the plant does not directly participate in. In the example that follows next, a similar but theoretical scenario between three agents is studied where two agents are initially interacting but not producing anything with a third agent on the sidelines. Then the third agent provides a signal to the other two to selectively change their behaviour that directly benefits two of the agents and indirectly benefits the third.

## 4.1 Model scenario

We begin with three agents in proximity to one another situated in a noisy environment in which it is possible for them to collectively do something useful but they are initially in the unfortunate situation in which the behaviour of the three agents produces nothing that is of value. The possible actions ($x_i$) of the agents ($A_i$) are binary: $x_i \in \{-1, 1\}$, $i \in \{1, 2, 3\}$ and so in this example the $x_i$ are not probabilities. $A_1$ is randomly and uniformly changing its states due to a (useful, information carrying) signal it receives from the environment at time $t$: $s^t \in \{-1, 1\}$ such that $P(s^t = 1) = P(s^t = -1) = 0.5$. $A_2$ and $A_3$ are initially engaged in the prisoner's dilemma (PD) game where, as a consequence of selfishly (naïvely) optimising their choice of $x_i$, they are in the defect-defect Nash equilibrium (NE), and so their joint action is constant (with zero value). They are capable of a second output when in the cooperate-cooperate configuration, and this output has a positive value. However, being stuck in the PD NE, they are not initially able to produce it.

$A_2$ and $A_3$ can produce something of value when they cooperate, but it only has value if the external signal $s^t = +1$. The output value of $A_2$ and $A_3$ interacting is initially a sequence of 0s: they are not cooperating with each other, and nothing of value is being produced. We assign each agent a state at time $t$, $x_i^t \in \{-1, 1\}$, such that their dynamic is an ordered sequence of binary states $[x_i^1, x_i^2, \ldots, x_i^T]$ for $t \in \{1, \ldots, T\}$. The output from $A_2$ and $A_3$ interacting at time $t$ is a result of having matched their states: $V(o^t) = 1$ if $x_2^t = x_3^t = 1$ and $V(o^t) = 0$ if $x_2^t = x_3^t = -1$, the $x_2^t \neq x_3^t$ cases are never achieved as they are not NE and $A_2$ and $A_3$ will always choose according to the NE of their interactions. The collective behavioural vector is: $\mathbf{x} = \{x_1, x_2, x_3\}$ or time indexed: $\mathbf{x}^t = \{x_1^t, x_2^t, x_3^t\}$.

The signal $s^t$ that $A_1$ receives indicates whether or not $A_2$ and $A_3$ should cooperate at time $t$, but initially no information is passing from $A_1$ to $\{A_2, A_3\}$ and $A_1$ cannot produce anything by itself. So $U_1(\mathbf{x}; \mathbf{a}) = 0$, and agents $A_2$ and $A_3$ only produce output of zero value: $V(o^t) = 0$. The agent utilities in this case are:

$$U_1(\mathbf{x}; \mathbf{a}) = 0, \qquad\qquad U_2(\mathbf{x}; \mathbf{a}) = a_2^0 + a_2^3 x_3 + a_2^2 x_2, \qquad\qquad U_3(\mathbf{x}; \mathbf{a}) = a_3^0 + a_3^2 x_2 + a_3^3 x_3, \tag{7}$$

and $U_2 = U_3 = 0$ for defect-defect joint strategies. This is equivalent to the configuration in Figure 1(a), and we illustrate this in Figure 2, a. In the second scenario, we assume that $A_2$ and $A_3$ are still naïvely pursuing their selfish goals, but
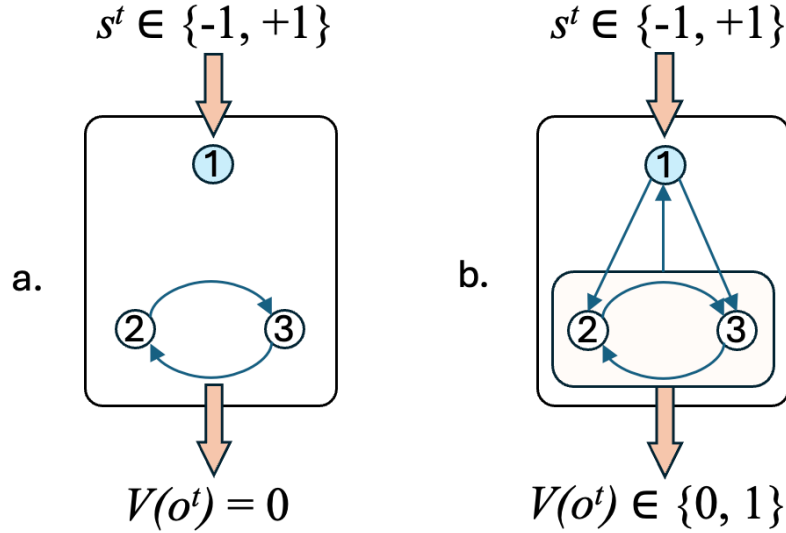
Figure 2: The interaction networks for the two scenarios: (a) $A_1$ perceives signal $s^t$ but cannot relay it to $A_2$ and $A_3$ and these two agents are not cooperating and so produce 0 value from their output $o^t$. (b) $A_1$ influences $A_2$ and $A_3$ to cooperate when $A_1$ receives a signal from the environment that cooperating will result in a positive payoff for $A_2$ and $A_3$; as a consequence, $A_1$ receives a portion of the utility from both $A_2$ and $A_3$.

|  |  | $A_3 \to x_3$ | |
|---|---|---|---|
|  |  | $+1 \equiv$ C | $-1 \equiv$ D |
| $A_2 \to x_2$ | $+1 \equiv$ Cooperate | (R, R) | (S, T) |
|  | $-1 \equiv$ Defect | (T, S) | (P, P) |

(a) Generalised payoff matrix

|  | C | D |
|---|---|---|
| C | $(1, 1)$ | $(0 - c, 1 + c)$ |
| D | $(1 + c, 0 - c)$ | $(0, 0)$ |

(b) Prisoner's Dilemma $(c = -\frac{1}{4})$ and Harmony $(c = \frac{1}{4})$

Figure 3: Payoff matrices for agents $A_2$ and $A_3$, which agent $A_1$ can strategically influence by setting $c \leftarrow x_1$.

now $A_1$ is more psychologically aware: it knows the beliefs, preferences, and constraints of the other two agents, and it is able to influence their preferences so that they will cooperate with each other when $A_1$ receives the signal: $s^t = 1$. They, in turn, will pay $A_1$ some portion of their total payoff.

To do this numerically, we standardise the symmetric, two-player, two-choice, game theory payoff matrices for $A_2$ and $A_3$ by setting $R = 1$ and $P = 0$ (Figure 3, a. and see [117]). We note that the Harmony game (Ha), in which both agents coordinating is the only Nash equilibrium and $T > R > P > S$, can be transformed into the Prisoner's Dilemma (PD), in which both agents defecting is the only Nash equilibrium and $R > T > S > P$, by introducing a two-state parameter $c \in \{-\frac{1}{4}, \frac{1}{4}\} \implies \{\text{PD}, \text{Ha}\}$ and replacing $T = R + c$ and $S = P - c$ so that the payoff matrices derivable from Figure 3,b. result in the utilities:

$$U_2(\mathbf{x}; \mathbf{a}) = R + cx_2 - (R + c)x_3, \qquad U_3(\mathbf{x}; \mathbf{a}) = R + cx_3 - (R + c)x_2. \qquad (8)$$

As agent $A_1$ is aware of the specific BPC model by which $A_2$ and $A_3$ make their decisions, $A_1$ adjusts their behaviour by manipulating their utility co-factors following $s^t$, by setting $x_1^t = \frac{1}{4}s^t$ such that $x_1^t \in \{-\frac{1}{4}, \frac{1}{4}\}$. $A_1$ then modulates the utility of the other two agents (via an information-carrying signal such as speaking to them) and, in response, receives 10% of the payoffs that result from the interaction between $A_2$ and $A_3$. The resulting payoffs for each agent are (derived in detail in Appendix B):

$$U_1(\mathbf{x}; \mathbf{a}) = \frac{1}{10}\big(U_2(\mathbf{x}; \mathbf{a}) + U_3(\mathbf{x}; \mathbf{a})\big), \quad U_2(\mathbf{x}; \mathbf{a}) = R + x_1x_2 - (R + x_1)x_3, \quad U_3(\mathbf{x}; \mathbf{a}) = R + x_1x_3 - (R + x_1)x_2. \quad (9)$$

## 4.2 Model interpretation

We make three observations regarding these two scenarios. First, we note that the interactions in the second scenario are a hypergraph in the sense that if any one of the nodes were removed in Figure 2, b. no remaining agent receives any payoff, all agents are necessary contributors to any single agent receiving a utility for their actions [115]. Second, it is straightforward to compute the collective intelligence of both scenarios using Equation 1 and the three stochastic variables $\{x_1^t, x_2^t, x_3^t\}$ when $s^t \in \{-1, 1\}$ is uniformly distributed: In the first scenario it is 0 bits and for the second

scenario it is 1 bit. This is a simple illustration of how a single agent, knowing the BPC model of two other agents, who are both naïve to all other agents' BPC models, can improve the collective intelligence of the group by manipulating the network of interactions between agents. This is analogous to the results discussed in Woolley [54] in which group performance, using a measure of collective intelligence, is improved by having agents who have a theory of mind. Third, unlike evolutionary processes, the process in the second scenario has a goal-directed agent who is psychologically informed about the internal states of the other agents and so is able to make plans contingent on different configurations of the interaction network. This implies that $A_1$ understands the preferences implied by the co-factors of Equation 8, $(R + c)$ and Equation 9, $(R + x_1)$, influencing the behaviour of the other two agents. These preferences are the hidden (cognitive) variables discussed in Section 2.2, i.e. the *Preferences* aspect of the BPC framework. In terms of levels of ToM, $A_1$ only needs to ascribe hidden states to the other two agents but not the ability for $A_2$ or $A_3$ to ascribe hidden states to $A_1$, for example the NE can be reached by $A_2$ and $A_3$ by observing the other agent's behaviour. So the ToM of $A_1$ is of second order and $A_2$ and $A_3$ only need a first order ToM. In this sense a ToM allows $A_1$ to reconfigure their social networks to further their goals at much shorter time scales than the evolutionary time scale of ecological networks.

# 5 Strategic Inference, Theory of Mind, and the Dual Role of `piKL`

In strategic multi-agent environments, coordination and collaboration often depend not only on optimal decision-making, but on the agent's ability to reason about the intentions, beliefs, and goals of others. This capability—often referred to as *Theory of Mind* (ToM)—becomes even more powerful when paired with language as a medium for influencing and inferring mental states. In this section, we develop a formal model that integrates reinforcement learning, ToM, and communication into a unified agent architecture. We also clarify two distinct interpretations of the `piKL` algorithm introduced in the Cicero AI framework.

## 5.1 Baseline Reinforcement Learning Policy

Let $s \in \mathcal{S}$ denote the environment state and $a \in \mathcal{A}$ the action selected by agent $i$. The baseline policy derived from reinforcement learning (RL) maximizes expected return:

$$\pi^{\text{RL}}(a \mid s) = \arg\max_a Q(s, a)$$

where $Q(s, a)$ is the value function learned via standard temporal-difference or policy-gradient methods.

## 5.2 Theory of Mind and Noisy Communication Channels

Each agent $i$ is assumed to have an internal, latent cognitive state $\theta_i \in \Theta$, encoding beliefs, goals, or preferences. Suppose agent $i$ communicates with agent $j$ via a message $m \in \mathcal{M}$ drawn from a noisy channel:

$$P(m \mid \theta_i)$$

The receiving agent $j$ updates its belief over $\theta_i$ using Bayes' rule:

$$b_j(\theta_i \mid m) = \frac{P(m \mid \theta_i) \cdot P(\theta_i)}{\sum_{\theta_i'} P(m \mid \theta_i') \cdot P(\theta_i')}$$

This belief over the sender's internal state allows agent $j$ to modulate its action policy accordingly.

## 5.3 ToM-Enhanced Policy via Belief Inference

The ToM-enhanced policy for agent $j$ is defined by taking the expectation over the inferred mental states of agent $i$:

$$\pi_j^{\text{ToM}}(a \mid s, m) = \mathbb{E}_{\theta_i \sim b_j(\theta_i \mid m)} [\pi_j(a \mid s, \theta_i)]$$

This formulation permits the agent to adaptively align, cooperate, or compete based on its belief about the intentions or preferences of the other agent.

## 5.4 Strategic Pragmatics and Message Selection

Recognizing that the receiver is performing this inference, the speaker agent $i$ may choose its message strategically to manipulate beliefs:

$$m^* = \arg\max_m \mathbb{E}_{\theta_j \sim b_i(\theta_j)} [U_i(\theta_j(m))]$$

This recursive, belief-aware messaging strategy is analogous to pragmatic language models such as the Rational Speech Acts (RSA) framework.

## 5.5 Two Interpretations of `piKL`

The `piKL` method measures divergence between policies. We describe two distinct formulations:

**1. Anchor-Regularised `piKL` (Cicero)** In the Cicero system, the agent's current policy $\pi_i$ is constrained by a divergence from an anchor policy $\pi_{\text{anchor}}$, learned from human gameplay:

$$D_{\text{piKL}}^{\text{anchor}} = \text{KL}\left[\pi_i(a \mid s) \,\|\, \pi_{\text{anchor}}(a \mid s)\right]$$

This encourages the agent to remain near human-like, interpretable behavior. The total objective becomes:

$$\mathcal{L}_{\text{anchor}} = \mathbb{E}_{(s,a)\sim\pi_i}\left[R(s,a)\right] - \lambda_1 \cdot D_{\text{piKL}}^{\text{anchor}}$$

**2. ToM-Induced `piKL`** Alternatively, we define a divergence between the baseline RL policy and the ToM-enhanced policy:

$$D_{\text{piKL}}^{\text{ToM}} = \text{KL}\left[\pi^{\text{RL}}(a \mid s) \,\|\, \pi^{\text{ToM}}(a \mid s, m)\right]$$

This measures how much the agent's behavior is changed by incorporating belief inference and communicative reasoning.

## 5.6 Unified Objective

We combine both interpretations in a unified loss:

$$\mathcal{L}_{\text{total}} = \mathbb{E}_{(s,a)\sim\pi_i}\left[R(s,a)\right] - \lambda_1 \cdot \text{KL}\left[\pi_i \,\|\, \pi_{\text{anchor}}\right] - \lambda_2 \cdot \text{KL}\left[\pi^{\text{RL}} \,\|\, \pi^{\text{ToM}}\right]$$

Here, $\lambda_1$ and $\lambda_2$ control the influence of human-alignment and social-inference regularisation, respectively. Together, they allow an agent to balance exploitation, interpretability, and strategic communicative inference.

## 5.7 Interpretation

This model formalizes a cognitive architecture in which:

- Reinforcement learning drives raw goal-seeking behavior;

- Theory of Mind allows adaptation to others' latent beliefs;

- Language serves as a noisy yet informative conduit of cognitive intent;

- Pragmatic reasoning aligns communicative goals with strategic planning;

- Dual regularisation (via `piKL`) ensures the agent balances optimality with trustworthiness and social intelligence.

## 5.8 Our perspective in context

We have reviewed some of the recent advances in how people, as complex agents, with a vast space of specialised, context-dependent behaviours, are able to coordinate their activities. In particular, finding ways to recombine our individual competencies to produce, in a short period of time, the appropriate collective competencies is a combinatorically complex task. The central theme of this article is that having a causal (generative [118, 119]) model of another person's behaviour that reflects the current environmental and social context helps us manage this complexity. This is one way in which our solid brains produce the necessary liquid social structures that quickly build collective solutions with novel collective competencies. In this section we bring these ideas together and summarise our view.

We first draw attention to how, in Section 4, $A_1$ has constructed a niche for itself in the context of the preexisting dyadic network between $A_2$ and $A_3$. In ecological networks a new agent joins a pre-existing network if it can find a niche within which it can fit. This occurs in one of three distinct ways: niche *choice*, niche *conformance*, and niche *construction* [120, 121]. Niche choice occurs when an individual selects environmental conditions that align with its phenotype, while niche conformance involves adjusting its phenotype to suit the environment. Niche construction is the modification of the environment to meet individual needs, which may also impact other species. This suggests an analogy between the formation of ecological networks, where new agents joining may enhance or disrupt the current configuration, and social networks where membership can be explicitly or implicitly gated according to a prospective agent's 'fit' within the group, and even if the group can adjust to accommodate a newcomer, just as the newcomer can adjust in order to be accepted. In human social groups people can be recruited or excluded depending on their contribution to the better functioning of the collective, which in turn corresponds to changes in individual neural activity related to the social network structure [122].

In our analogy with ecological networks, we suggest that *any* signalling agents in a network need to encode messages over an information-carrying medium that receivers are receptive to and that can then be decoded by a specific receiver. We have argued that humans, as both signallers and receivers, take advantage of our ToM in order to understand what signals will be correctly interpreted by another person and to adapt their signalling to the cognitive state of the receiver. Specifically, our psychology allows us to learn how to read, interpret, and re-write our interpersonal communication over a very short time frame, allowing us rapid, targeted control over our collective behaviour.

In some sense the adage *there is nothing new under the sun* holds here as evolution and biology have recycled fundamental, pre-existing principles in the service of human sociality. But the adaptive speed of our social networks, the psychological mechanisms involved, the variety of purposes they serve, and the complexity of the communications appear most highly developed in humans. In ecological networks, for example, some agents can act as an encoder-sender of semantic information via infochemicals that signal another agent. A second agent then acts as a receiver-decoder of this signal that in turn changes the behaviour of the receptive agent. The combination of the receiver's anatomical configuration and behavioural phenotype elicits an appropriate stimulus-response that benefits the signalling agent, and these interactions can form vast, complex hypergraphs of competition and cooperation between agents of many different species. But neither agent needs higher cognitive faculties—their ToM is of the zeroth order. By analogy, humans can target specific individuals or groups of individuals in order to have them adjust their behaviours to best suit the goals of the signaller. The competencies that a ToM affords the sender allows them to know that a receiver is capable of both decoding the signal and acting appropriately in response because they know of the receiver's *receptive* and *causal* states, both cognitive (hidden) and behavioural (overt). That is to say they take advantage of their ToM to understand how, when, what, and to whom they need to signal in order to achieve a beneficial outcome, whereas blind evolution would take much longer to achieve the same results.

The separation of time scales also plays an important distinction between learning processes that use the same theoretical foundation. The mathematical description of the processes that underpin the evolution of species has been shown to be equivalent to the statistical learning process that underpins Bayesian learning in individual agents [123, 124, 125], for both biological and artificial agents. What Bayesian learning provides at the level of the individual agent is an advantage in the speed of adaptation that Bayesian learning via evolution cannot achieve. Again, this suggests a universality of description that only varies qualitatively in terms of time scales and the specifics of the mechanisms.

The models we presented highlight the roles ToM plays in human social collectives such as deciding which relationships to develop and how these relationships aid in the collective processing of information towards an end goal. Being able to understand, predict, and manipulate these information-carrying connections is valuable in improving group dynamics (in human groups, AI groups, or hybrid human-AI groups) and for the future design of ToM-based AI. For example, knowing that ToM plays a role in group performance [54], quantifying and modelling (ToM-mediated) CI will shed light on the types of interventions or adjustments that can modulate it. Additionally, including AI agents with sophisticated ToM in human groups has the potential to 'construct social niches' (see below) more effectively, such as by modifying the incentives involved in a social setting (making some actions more desirable than others; cf. mechanism design [126]), or the connections between agents [127], or by simply mediating interactions with an increased capacity for memory, attention, and reasoning [128]. Recent efforts have leveraged ToM and AI for improving team effectiveness [129, 128] and understanding the link between ToM and deception [130], or ToM and expectation formation in markets [131]. Others suggest AI with ToM could assist humans in improving their ToM and communications skills [132], in negotiation, education (adaptation to students needs and knowledge), and games [133]; healthcare (to tailor care to patients' mental states), self-driving cars (anticipating the actions of other cars), workplaces (employee mood and stress), and marketing [134].

To successfully apply ToM-based AI to these areas, we need to consider how we could define the appropriate incentives for an AI agent to adapt to a given set of social circumstances in a desirable manner. Our examples, while minimal, offer a conceptual framing for thinking about this human-AI integration. For example, the use of information theory to quantify the information processing provides a measure that can be used in an optimisation objective for an AI agent within a collective, incentivising an AI to act in ways that are constructive for the collective as a whole.

The difficulties of implementing socially aware algorithms in an AI are also made more complicated by human variations in expectations. For example, the famous trolley problem has often been used as an example of a moral dilemma that an autonomous AI that is in control of a car may confront during an emergency. To understand the complexities of this issue for an AI, Awad et al [135] collected 40 million individual decisions from 233 countries and territories in 10 languages. The questions focused on trolley-like moral dilemmas to examine cross-cultural ethical variation, showing three major clusters of countries and that these variations reflect differences in modern institutions and deep cultural traits. These results illustrate that our morality is not universal, there is a rich variety of expectations for the same moral dilemmas, and that any AI would need to reflect the local cultural variations in which they are used. In the Discussion next we cover the broader issue of how an AI will need to learn and adapt in order to fit the social circumstances the AI will be expected to be constrained by, very much like other socio-ecological systems.

To date the study of the collective intelligence in hybrid AI-human systems is in its infancy, however we already have many of the necessary tools to begin these investigations. ToM is experimentally well established in psychology and

work has begun developing and implementing algorithms for an AI-ToM, see for example Mao et al's recent review of empirical work on beliefs, desires and preferences for AI with ToM [136] as well as the computational, algorithmic, and experimental articles discussed in the Introduction above. What is not as well understood are the dynamics that agential AI will introduce into our human social ecology, and to that extent the study of *machine behaviour* [36] and new measures of collective intelligence as we have described here need to be applied to AI-human hybrid systems to better understand the collective outcomes, rather than focusing on the algorithmic foundations.

# 6  Discussion: Theory of Mind in the Context of Social AI

Each step in our cultural development has changed the ways in which we combine individual skills to achieve better, more sophisticated collective outcomes. There is evidence that hunter gatherers participated in the division of labour, complex social networks, multilevel, and fractal-like social structures long before we settled into villages [137, 138, 6, 4]. This was in part due to our ability to construct our own niches [139] and this practice has persisted from hunter-gathers to farming [140] through to civilisation building. As Arroyo-Kalin *et al.* [141] quote in their opening to the special issue *Civilisation and Human Niche Construction*:

> It is impossible to avoid the conclusion that organisms construct every aspect of their environment themselves. They are not the passive objects of external forces, but the creators and modulators of these forces. The metaphor of adaptation must therefore be replaced by one of construction, a metaphor that has implications for the form of evolutionary theory (Levins and Lewontin 1985: 104).

It has also been argued that through the manipulation of our environmental niche we have brought about the Anthropocene [142, 143, 144, 145].

On the other hand, *social niche construction* [146] extends this idea to the process whereby agents modify their social context so as to influence their own social evolution. A *social niche* is the context in which social behaviour occurs, and so social niche construction is where agents actively choose to change their social environments, for example, by choosing who to associate with and how to behave [139]. Ryan *et al.* [146] describe this in game theoretical terms as the *effective game* agents are playing after all relevant factors are accounted for, such as the underlying game itself (the payoff matrices in conventional games for example) and any *social niche modifiers*, amongst others. A *social niche modifier* is a trait that alters the effective game being played, causing it to differ from the immediate payoff matrices that are usually the complete description of the incentivised interactions, for example population structure, relatedness, punishment etc. In the example in Section 4, $A_1$ constructs a social niche for itself by manipulating the structure of the game that $A_2$ and $A_3$ are playing, for the benefit of $A_1$ and incidentally for the benefit of the other agents as well. However, social niche construction theory pertains to evolution more broadly and is not specific to human social networks as the formal description of the model in Section 4 could be an evolved network of plants and animals or a more rapidly changing social network.

In moving towards AI that is situated within a hybrid human-AI ecology, the complexities of effective network construction, communication, and cognitive tools will need to be worked through. Here, we discuss just two of the issues: the difficulties of building psychologically complex AIs and the social environment AI will need to adapt to. The capacity of current AI theories to be sufficient to encompass, in principle at least, human levels of cognition is a rich area of research, from attention mechanisms [147], to reinforcement learning subsuming reasoning [148] and chain of thought reasoning in language models [149] there are arguments being made for *emergent phenomena* in AI [150] as well as ToM [151] and even artificial general intelligence [152]. What is often missed in these examples is that the individual AI's ability to digest information and update parameter weights is only a small fraction of what is needed to be effective in a specifically social context. An important tool in our psychological toolbox is our ability to maintain shared hidden variables that are the causal basis of our coordinated joint actions in physical and social environments. Even progress on single agents having effective causal models of the physical environment has been much slower than in other areas of AI [153]. As we have shown in this article, there is very good evidence that causal models, both physical and social, are necessary for people to be able to communicate with each other, and communication is inextricably tied to our ToM. This triad of language, shared causal models of hidden variables, and ToM appears to be a minimum for human social coordination.

Even with this triad established within an AI, the next challenge is where, when, and how an AI should fit into any given collaborative social context. Simple machine intelligence in collaboration with people, such as auto-correct, recommender systems, or GPS navigation, are effective because a human decided there was a need for these tools, they placed them in the appropriate context and then the users adjust their behaviour to any shortcomings the tool might have. But the more autonomous machine intelligence becomes, the more complex the physical and social environment is, and the more trust that needs to be placed in the causal models that drive the behaviour of an AI (i.e. their analogue of a person's hidden cognitive variables) the more an AI needs to know how to interact with us in a way that resembles how we interact with each other. They will ultimately need to be able to do this via niche construction, niche adaptation, and niche choice, all of which, for people, is a negotiated relationship between each other that needs to be satisficing for those involved. This brings human-AI joint adaptation closer to what Laland *et al.* have in mind when considering a rethink of evolutionary theory [154]:

We hold that organisms are constructed in development, not simply 'programmed' to develop by genes. Living things do not evolve to fit into pre-existing environments, but co-construct and coevolve with their environments, in the process changing the structure of ecosystems.

Consequently, we argue that for any AI to be a *beneficially adaptive*, autonomous, and socially aware agent, it will also need to reflect the universal principles we see in many evolutionary transitions [155, 156, 157]. This includes multi-level selection in economic [145, 158] and biological systems [159, 11] as well as the hierarchical transitions [159] that appear to be ubiquitous in our biological, social, economic, and technological history.

# References

[1] Michael Levin. Bioelectric networks: the cognitive glue enabling evolutionary scaling from physiology to mind. *Animal Cognition*, 26(6):1865–1891, 2023.

[2] Ida Momennejad. Collective minds: social network topology shapes collective cognition. *Philosophical Transactions of the Royal Society B*, 377(1843):20200315, 2022.

[3] Andrea Bamberg Migliano and Lucio Vinicius. The origins of human cumulative culture: from the foraging niche to collective intelligence. *Philosophical Transactions of the Royal Society B*, 377(1843):20200317, 2022.

[4] Marcus J Hamilton, Bruce T Milne, Robert S Walker, Oskar Burger, and James H Brown. The complex structure of hunter–gatherer social networks. *Proceedings of the Royal Society B: Biological Sciences*, 274(1622):2195–2203, 2007.

[5] Russell A Hill, R Alexander Bentley, and Robin IM Dunbar. Network scaling reveals consistent fractal pattern in hierarchical mammalian societies. *Biology letters*, 4(6):748–751, 2008.

[6] Michael S Harré and Mikhail Prokopenko. The social brain: scale-invariant layering of erdős–rényi networks in small-scale human societies. *Journal of the Royal Society Interface*, 13(118):20160044, 2016.

[7] Ricard Solé, Melanie Moses, and Stephanie Forrest. Liquid brains, solid brains, 2019.

[8] Jordi Piñero and Ricard Solé. Statistical physics of liquid brains. *Philosophical Transactions of the Royal Society B*, 374(1774):20180376, 2019.

[9] Michael Levin. The computational boundary of a "self": developmental bioelectricity drives multicellularity and scale-free cognition. *Frontiers in psychology*, 10:2688, 2019.

[10] Michael Levin. Technological approach to mind everywhere: an experimentally-grounded framework for understanding diverse bodies and minds. *Frontiers in systems neuroscience*, 16:768201, 2022.

[11] Michael Levin. Darwin's agential materials: evolutionary implications of multiscale competency in developmental biology. *Cellular and Molecular Life Sciences*, 80(6):142, 2023.

[12] Chris Fields and Michael Levin. Competency in navigating arbitrary spaces as an invariant for analyzing cognition in diverse embodiments. *Entropy*, 24(6):819, 2022.

[13] John J Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the national academy of sciences*, 79(8):2554–2558, 1982.

[14] David Beniaguev, Idan Segev, and Michael London. Single cortical neurons as deep artificial neural networks. *Neuron*, 109(17):2727–2739, 2021.

[15] Sophie Deneve. Bayesian spiking neurons i: inference. *Neural computation*, 20(1):91–117, 2008.

[16] Tomas Kay, Alba Motes-Rodrigo, Arthur Royston, Thomas O Richardson, Nathalie Stroeymeyt, and Laurent Keller. Ant social network structure is highly conserved across species. *Proceedings B*, 291(2027):20240898, 2024.

[17] Thomas O Richardson, Andrea Coti, Nathalie Stroeymeyt, and Laurent Keller. Leadership–not followership–determines performance in ant teams. *Communications biology*, 4(1):535, 2021.

[18] Sebastian Stockmaier, Nathalie Stroeymeyt, Eric C Shattuck, Dana M Hawley, Lauren Ancel Meyers, and Daniel I Bolnick. Infectious diseases and social distancing in nature. *Science*, 371(6533):eabc8881, 2021.

[19] Mirta Galesic, Daniel Barkoczi, Andrew M Berdahl, Dora Biro, Giuseppe Carbone, Ilaria Giannoccaro, Robert L Goldstone, Cleotilde Gonzalez, Anne Kandler, Albert B Kao, et al. Beyond collective intelligence: Collective adaptation. *Journal of the Royal Society interface*, 20(200):20220736, 2023.

[20] Philip Mirowski and Koye Somefun. Markets as evolving computational entities. *Journal of Evolutionary Economics*, 8(4):329–356, 1998.

[21] Robert L. Axtell. Economics as distributed computation. In *Meeting the Challenge of Social Problems via Agent-Based Simulation: Post-Proceedings of the Second International Workshop on Agent-Based Approaches in Economic and Social Complex Systems*, pages 3–23. Springer, 2003.

[22] Philip Mirowski. Markets come to bits: Evolution, computation and markomata in economic science. *Journal of Economic Behavior & Organization*, 63(2):209–242, 2007.

[23] Michael S Harré. Entropy, economics, and criticality. *Entropy*, 24(2):210, 2022.

[24] Michael S. Harré and Terrence Bossomaier. Phase-transition–like behaviour of information measures in financial markets. *Europhysics Letters*, 87(1):18009, 2009.

[25] Michael Harré. Entropy and transfer entropy: the dow jones and the build up to the 1997 asian crisis. In *Proceedings of the International Conference on Social Modeling and Simulation, plus Econophysics Colloquium 2014*, pages 15–25. Springer International Publishing, 2015.

[26] Chris G Langton. Computation at the edge of chaos: Phase transitions and emergent computation. *Physica D: nonlinear phenomena*, 42(1-3):12–37, 1990.

[27] Douglas G Moore, Gabriele Valentini, Sara I Walker, and Michael Levin. Inform: efficient information-theoretic analysis of collective behaviors. *Frontiers in Robotics and AI*, 5:60, 2018.

[28] Michael Wibral, Raul Vicente, and Joseph T Lizier. *Directed information measures in neuroscience*, volume 724. Springer, 2014.

[29] Giulio Tononi, Melanie Boly, Marcello Massimini, and Christof Koch. Integrated information theory: from consciousness to its physical substrate. *Nature reviews neuroscience*, 17(7):450–461, 2016.

[30] Pedro AM Mediano, Fernando E Rosas, Juan Carlos Farah, Murray Shanahan, Daniel Bor, and Adam B Barrett. Integrated information as a common signature of dynamical and information-processing complexity. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 32(1), 2022.

[31] Adam B Barrett and Anil K Seth. Practical measures of integrated information for time-series data. *PLoS computational biology*, 7(1):e1001052, 2011.

[32] Pedro AM Mediano, Anil K Seth, and Adam B Barrett. Measuring integrated information: Comparison of candidate measures in theory and simulation. *Entropy*, 21(1):17, 2018.

[33] Mikhail Prokopenko, Fabio Boschetti, and Alex J Ryan. An information-theoretic primer on complexity, self-organization, and emergence. *Complexity*, 15(1):11–28, 2009.

[34] Joseph T Lizier. *The local information dynamics of distributed computation in complex systems*. Springer Science & Business Media, 2012.

[35] Terry Bossomaier, Lionel Barnett, Michael Harré, and Joseph T Lizier. *An Introduction to Transfer Entropy: Information Flow in Complex Systems*. Springer, 2016.

[36] Iyad Rahwan, Manuel Cebrian, Nick Obradovich, Josh Bongard, Jean-François Bonnefon, Cynthia Breazeal, Jacob W Crandall, Nicholas A Christakis, Iain D Couzin, Matthew O Jackson, et al. Machine behaviour. *Nature*, 568(7753):477–486, 2019.

[37] Neil Johnson, Guannan Zhao, Eric Hunsader, Hong Qi, Nicholas Johnson, Jing Meng, and Brian Tivnan. Abrupt rise of new machine ecology beyond human response time. *Scientific reports*, 3(1):2627, 2013.

[38] Nadav Aharony, Wei Pan, Cory Ip, Inas Khayal, and Alex Pentland. Social fmri: Investigating and shaping social mechanisms in the real world. *Pervasive and mobile computing*, 7(6):643–659, 2011.

[39] Niko Tinbergen. On aims and methods of ethology. *Animal Biology*, 55(4), 2005.

[40] Jason W Burton, Ezequiel Lopez-Lopez, Shahar Hechtlinger, Zoe Rahwan, Samuel Aeschbach, Michiel A Bakker, Joshua A Becker, Aleks Berditchevskaia, Julian Berger, Levin Brinkmann, et al. How large language models can reshape collective intelligence. *Nature human behaviour*, pages 1–13, 2024.

[41] Jacob W Crandall, Mayada Oudah, Tennom, Fatimah Ishowo-Oloko, Sherief Abdallah, Jean-François Bonnefon, Manuel Cebrian, Azim Shariff, Michael A Goodrich, and Iyad Rahwan. Cooperating with machines. *Nature communications*, 9(1):233, 2018.

[42] David Sally. Conversation and cooperation in social dilemmas: A meta-analysis of experiments from 1958 to 1992. *Rationality and society*, 7(1):58–92, 1995.

[43] Gary A Klein. A recognition-primed decision (rpd) model of rapid decision making. *Decision making in action: Models and methods*, 5(4):138–147, 1993.

[44] Meta Fundamental AI Research Diplomacy Team (FAIR)†, Anton Bakhtin, Noam Brown, Emily Dinan, Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff, Jonathan Gray, Hengyuan Hu, et al. Human-level play in the game of diplomacy by combining language models with strategic reasoning. *Science*, 378(6624):1067–1074, 2022.

[45] Edd Gent. Machine mind readers. *New Scientist*, 257(3426):46–49, 2023.

[46] Gerd Gigerenzer. Psychological ai: Designing algorithms informed by human psychology. *Perspectives on Psychological Science*, 19(5):839–848, 2024.

[47] Hao Cui and Taha Yasseri. Ai-enhanced collective intelligence. *Patterns*, 5(11), 2024.

[48] Milena Tsvetkova, Taha Yasseri, Niccolo Pescetelli, and Tobias Werner. A new sociology of humans and machines. *Nature Human Behaviour*, 8(10):1864–1876, 2024.

[49] Maxinder S Kanwal, Joshua A Grochow, and Nihat Ay. Comparing information-theoretic measures of complexity in boltzmann machines. *Entropy*, 19(7):310, 2017.

[50] Richard Watson and Michael Levin. The collective intelligence of evolution and development. *Collective Intelligence*, 2(2):26339137231168355, 2023.

[51] Uta Frith and Chris Frith. The social brain: allowing humans to boldly go where no other species has been. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 365(1537):165–176, 2010.

[52] Derek C Penn and Daniel J Povinelli. On the lack of evidence that non-human animals possess anything remotely resembling a 'theory of mind'. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 362(1480):731–744, 2007.

[53] Christopher Krupenye and Josep Call. Theory of mind in animals: Current and future directions. *Wiley Interdisciplinary Reviews: Cognitive Science*, 10(6):e1503, 2019.

[54] Anita Williams Woolley, Christopher F Chabris, Alex Pentland, Nada Hashmi, and Thomas W Malone. Evidence for a collective intelligence factor in the performance of human groups. *science*, 330(6004):686–688, 2010.

[55] David Engel, Anita Williams Woolley, Lisa X Jing, Christopher F Chabris, and Thomas W Malone. Reading the mind in the eyes or reading between the lines? theory of mind predicts collective intelligence equally well online and face-to-face. *PloS one*, 9(12):e115212, 2014.

[56] Wako Yoshida, Ray J Dolan, and Karl J Friston. Game theory of mind. *PLoS computational biology*, 4(12):e1000254, 2008.

[57] Herbert Gintis. The foundations of behavior: The beliefs, preferences, and constraints model. *Biological Theory*, 1:123–127, 2006.

[58] Herbert Gintis. A framework for the unification of the behavioral sciences. *Behavioral and brain sciences*, 30(1):1–16, 2007.

[59] Herbert Gintis. Unifying the behavioral sciences ii. *Behavioral and brain sciences*, 30(1):45–53, 2007.

[60] Dino Pedreschi, Luca Pappalardo, Emanuele Ferragina, Ricardo Baeza-Yates, Albert-László Barabási, Frank Dignum, Virginia Dignum, Tina Eliassi-Rad, Fosca Giannotti, János Kertész, et al. Human-ai coevolution. *Artificial Intelligence*, page 104244, 2024.

[61] Samuel PL Veissière, Axel Constant, Maxwell JD Ramstead, Karl J Friston, and Laurence J Kirmayer. Thinking through other minds: A variational approach to cognition and culture. *Behavioral and brain sciences*, 43:e90, 2020.

[62] Karl J Friston, Maxwell JD Ramstead, Alex B Kiefer, Alexander Tschantz, Christopher L Buckley, Mahault Albarracin, Riddhi J Pitliya, Conor Heins, Brennan Klein, Beren Millidge, et al. Designing ecosystems of intelligence from first principles. *Collective Intelligence*, 3(1):26339137231222481, 2024.

[63] Michael S Harré and Husam El-Tarifi. Testing game theory of mind models for artificial intelligence. *Games*, 15(1):1, 2023.

[64] Michael S Harré. What can game theory tell us about an AI 'Theory of Mind'? *Games*, 13(3):46, 2022.

[65] Rachel Wood, Alexander G Liu, Frederick Bowyer, Philip R Wilby, Frances S Dunn, Charlotte G Kenchington, Jennifer F Hoyal Cuthill, Emily G Mitchell, and Amelia Penny. Integrated records of environmental change and evolution challenge the cambrian explosion. *Nature ecology and evolution*, 3(4):528–538, 2019.

[66] František Baluška and Michael Levin. On having no head: Cognition throughout biological systems. 7, 2016.

[67] A Boussard, J Delescluse, A Pérez-Escudero, and A Dussutour. Memory inception and preservation in slime moulds: the quest for a common mechanism. *Philosophical transactions of the Royal Society of London. Series B. Biological sciences*, 374(1774):20180368–20180368, 2019.

[68] Richard F. Thompson. Habituation: A history. *Neurobiology of learning and memory*, 92(2):127–134, 2009.

[69] Zhi Li and Satish K. Nair. Quorum sensing: How bacteria can coordinate activity and synchronize their response to external signals? *Protein science*, 21(10):1403–1417, 2012.

[70] Andrea Avena-Koenigsberger, Joaquín Goñi, Ricard Solé, and Olaf Sporns. Network morphospace. *Journal of the Royal Society Interface*, 12(103):20140881, 2015.

[71] Thomas Claverie and Peter C Wainwright. A morphospace for reef fishes: elongation is the dominant axis of body shape evolution. *PloS one*, 9(11):e112732, 2014.

[72] Luís F Seoane and Ricard Solé. The morphospace of language networks. *Scientific reports*, 8(1):10465, 2018.

[73] Michael Levin. Collective intelligence of morphogenesis as a teleonomic process. 2022.

[74] Richard A. Watson, Rob Mills, and C. L. Buckley. Global Adaptation in Networks of Selfish Components: Emergent Associative Memory at the System Scale. *Artificial Life*, 17(3):147–166, July 2011.

[75] Richard A. Watson, C. L. Buckley, and Rob Mills. Optimization in "self-modeling" complex adaptive systems. *Complexity*, 16(5):17–26, 2011.

[76] Markus Brede and Guillermo Romero-Moreno. Sensing enhancement on social networks: The role of network topology. *Entropy*, 24(5):738, 2022.

[77] Albert B Kao and Iain D Couzin. Modular structure within groups causes information loss but can improve decision accuracy. *Philosophical transactions of the Royal Society of London. Series B. Biological sciences*, 374(1774):20180378–20180378, 2019.

[78] Graham E Budd. Morphospace. *Current Biology*, 31(19):R1181–R1185, 2021.

[79] Xerxes D Arsiwalla, Ricard Sole, Clement Moulin-Frier, Ivan Herreros, Marti Sanchez-Fibla, and Paul Verschure. The morphospace of consciousness. *arXiv preprint arXiv:1705.11190*, 2017.

[80] Yasaman Bahri, Jonathan Kadmon, Jeffrey Pennington, Sam S Schoenholz, Jascha Sohl-Dickstein, and Surya Ganguli. Statistical mechanics of deep learning. *Annual Review of Condensed Matter Physics*, 11(1):501–528, 2020.

[81] Aina Ollé-Vila, Salva Duran-Nebreda, Núria Conde-Pueyo, Raúl Montañez, and Ricard Solé. A morphospace for synthetic organs and organoids: the possible and the actual. *Integrative Biology*, 8(4):485–503, 2016.

[82] Joshua C Peterson, David D Bourgin, Mayank Agrawal, Daniel Reichman, and Thomas L Griffiths. Using large-scale experiments and machine learning to discover theories of human decision-making. *Science*, 372(6547):1209–1214, 2021.

[83] Julian Jara-Ettinger. Theory of mind as inverse reinforcement learning. *Current Opinion in Behavioral Sciences*, 29:105–110, 2019.

[84] Jaime Ruiz-Serra and Michael S Harré. Inverse reinforcement learning as the algorithmic basis for theory of mind: current methods and open problems. *Algorithms*, 16(2):68, 2023.

[85] Thuy Ngoc Nguyen and Cleotilde Gonzalez. Theory of mind from observation in cognitive models and humans. *Topics in cognitive science*, 14(4):665–686, 2022.

[86] Michelle Zhao, Fade R Eadeh, Thuy-Ngoc Nguyen, Pranav Gupta, Henny Admoni, Cleotilde Gonzalez, and Anita Williams Woolley. Teaching agents to understand teamwork: Evaluating and predicting collective intelligence as a latent variable via hidden markov models. *Computers in Human Behavior*, 139:107524, 2023.

[87] Pranav Gupta, Thuy Ngoc Nguyen, Cleotilde Gonzalez, and Anita Williams Woolley. Fostering collective intelligence in human–ai collaboration: laying the groundwork for cohumain. *Topics in cognitive science*, 2023.

[88] Garriy Shteynberg, Jacob B Hirsh, Wouter Wolf, John A Bargh, Erica J Boothby, Andrew M Colman, Gerald Echterhoff, and Maya Rossignac-Milon. Theory of collective mind. *Trends in Cognitive Sciences*, 27(11):1019–1031, 2023.

[89] Michael Shum, Max Kleiman-Weiner, Michael L Littman, and Joshua B Tenenbaum. Theory of minds: Understanding behavior in groups through inverse planning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 6163–6170, 2019.

[90] Marlize Lombard and Peter Gärdenfors. Causal cognition and theory of mind in evolutionary cognitive archaeology. *Biological Theory*, 18(4):234–252, 2023.

[91] Miriam Noël Haidle. Working-memory capacity and the evolution of modern cognitive potential: implications from animal and early human tool use. *Current anthropology*, 51(S1):S149–S166, 2010.

[92] Jörg Lang, Jutta Winsemann, Dominik Steinmetz, Ulrich Polom, Lukas Pollok, Utz Böhner, Jordi Serangeli, Christian Brandes, Andrea Hampel, and Stefan Winghart. The pleistocene of schöningen, germany: a complex tunnel valley fill revealed from 3d subsurface modelling and shear wave seismics. *Quaternary Science Reviews*, 39:86–105, 2012.

[93] Francesca Happé. Theory of mind and the self. *Annals of the New York Academy of Sciences*, 1001(1):134–144, 2003.

[94] Jonathan Birch, Alexandra K Schnell, and Nicola S Clayton. Dimensions of animal consciousness. *Trends in cognitive sciences*, 24(10):789–801, 2020.

[95] Noam Chomsky. Biolinguistic explorations: Design, development, evolution. *International Journal of Philosophical Studies*, 15(1):1–21, 2007.

[96] M Jeffrey Farrar and Lisa Maag. Early language development and the emergence of a theory of mind. *First language*, 22(2):197–213, 2002.

[97] Jill G de Villiers. The role (s) of language in theory of mind. In *The neural basis of mentalizing*, pages 423–448. Springer, 2021.

[98] Jill G De Villiers and Jennie E Pyers. Complements to cognition: A longitudinal study of the relationship between complex syntax and false-belief-understanding. *Cognitive development*, 17(1):1037–1060, 2002.

[99] Karen Milligan, Janet Wilde Astington, and Lisa Ain Dack. Language and theory of mind: Meta-analysis of the relation between language ability and false-belief understanding. *Child development*, 78(2):622–646, 2007.

[100] Paul H Thibodeau and Lera Boroditsky. Metaphors we think with: The role of metaphor in reasoning. *PloS one*, 6(2):e16782, 2011.

[101] Paul H Thibodeau, Rose K Hendricks, and Lera Boroditsky. How linguistic metaphor scaffolds reasoning. *Trends in cognitive sciences*, 21(11):852–863, 2017.

[102] Eytan Bakshy, Itamar Rosenn, Cameron Marlow, and Lada Adamic. The role of social networks in information diffusion. In *Proceedings of the 21st international conference on World Wide Web*, pages 519–528, 2012.

[103] Robin IM Dunbar. Coevolution of neocortical size, group size and language in humans. *Behavioral and brain sciences*, 16(4):681–694, 1993.

[104] Robin IM Dunbar and Susanne Shultz. Evolution in the social brain. *science*, 317(5843):1344–1347, 2007.

[105] Robin IM Dunbar and Matt Spoors. Social networks, support cliques, and kinship. *Human nature*, 6:273–290, 1995.

[106] Tatiana Lau, Hillard T Pouncy, Samuel J Gershman, and Mina Cikara. Discovering social groups via latent structure learning. *Journal of Experimental Psychology: General*, 147(12):1881, 2018.

[107] Simon Baron-Cohen, Therese Jolliffe, Catherine Mortimore, and Mary Robertson. Another advanced test of theory of mind: Evidence from very high functioning adults with autism or asperger syndrome. *Journal of Child psychology and Psychiatry*, 38(7):813–822, 1997.

[108] Simon Baron-Cohen, Sally Wheelwright, Jacqueline Hill, Yogini Raste, and Ian Plumb. The "reading the mind in the eyes" test revised version: a study with normal adults, and adults with asperger syndrome or high-functioning autism. *The Journal of Child Psychology and Psychiatry and Allied Disciplines*, 42(2):241–251, 2001.

[109] Michael S Harré. Strategic information processing from behavioural data in iterated games. *Entropy*, 20(1):27, 2018.

[110] Daeyeol Lee, Michelle L Conroy, Benjamin P McGreevy, and Dominic J Barraclough. Reinforcement learning and decision making in monkeys during a competitive game. *Cognitive brain research*, 22(1):45–58, 2004.

[111] Luís MA Bettencourt. The rules of information aggregation and emergence of collective intelligent behavior. *Topics in Cognitive Science*, 1(4):598–620, 2009.

[112] Joseph T Lizier. Jidt: An information-theoretic toolkit for studying the dynamics of complex systems. *Frontiers in Robotics and AI*, 1:11, 2014.

[113] Jianing Wei, Lizhong Wang, Junwei Zhu, Sufang Zhang, Owi I Nandi, and Le Kang. Plants attract parasitic wasps to defend themselves against insect pests by releasing hexenol. *PLOS one*, 2(9):e852, 2007.

[114] Yutaka Kobayashi and Norio Yamamura. Evolution of signal emission by uninfested plants to help nearby infested relatives. *Evolutionary Ecology*, 21:281–294, 2007.

[115] Antonio J Golubski, Erik E Westlund, John Vandermeer, and Mercedes Pascual. Ecological networks over the edge: hypergraph Trait-Mediated Indirect Interaction (TMII) structure. *Trends in ecology & evolution*, 31(5):344–354, 2016.

[116] Adam Harris, Scott McCallum, and Michael S Harré. On the smooth unfolding of bifurcations in quantal-response equilibria. *Games and Economic Behavior*, 2023.

[117] Michael S Harré. Multi-agent economics and the emergence of critical markets. *arXiv preprint arXiv:1809.01332*, 2018.

[118] Karl Friston, Lancelot Da Costa, Noor Sajid, Conor Heins, Kai Ueltzhöffer, Grigorios A Pavliotis, and Thomas Parr. The free energy principle made simpler but not too simple. *Physics Reports*, 1024:1–29, 2023.

[119] Jaime Ruiz-Serra, Patrick Sweeney, and Michael S Harré. Factorised active inference for strategic multi-agent interactions. *arXiv preprint arXiv:2411.07362*, 2024.

[120] Andrew D Clark, Dominik Deffner, Kevin Laland, John Odling-Smee, and John Endler. Niche construction affects the variability and strength of natural selection. *The American Naturalist*, 195(1):16–30, 2020.

[121] Rose Trappes, Behzad Nematipour, Marie I Kaiser, Ulrich Krohs, Koen J Van Benthem, Ulrich R Ernst, Jürgen Gadau, Peter Korsten, Joachim Kurtz, Holger Schielzeth, et al. How individualized niches arise: Defining mechanisms of niche construction, niche choice, and niche conformance. *BioScience*, 72(6):538–548, 2022.

[122] Ralf Schmälzle, Matthew Brook O'Donnell, Javier O Garcia, Christopher N Cascio, Joseph Bayer, Danielle S Bassett, Jean M Vettel, and Emily B Falk. Brain connectivity dynamics during social interaction reflect social network structure. *Proceedings of the National Academy of Sciences*, 114(20):5153–5158, 2017.

[123] John O Campbell. Universal darwinism as a process of bayesian inference. *Frontiers in systems neuroscience*, 10:49, 2016.

[124] Richard A Watson, Rob Mills, CL Buckley, Kostas Kouvaris, Adam Jackson, Simon T Powers, Chris Cox, Simon Tudge, Adam Davies, Loizos Kounios, et al. Evolutionary connectionism: algorithmic principles underlying the evolution of biological organisation in evo-devo, evo-eco and evolutionary transitions. *Evolutionary biology*, 43:553–581, 2016.

[125] Richard A Watson and Eörs Szathmáry. How can evolution learn? *Trends in ecology & evolution*, 31(2):147–157, 2016.

[126] Steve Phelps, Peter McBurney, and Simon Parsons. Evolutionary mechanism design: a review. *Autonomous agents and multi-agent systems*, 21(2):237–264, 2010.

[127] Kevin R. McKee, Andrea Tacchetti, Michiel A. Bakker, Jan Balaguer, Lucy Campbell-Gillingham, Richard Everett, and Matthew Botvinick. Scaffolding cooperation in human groups with deep reinforcement learning. *Nature Human Behaviour*, 7(10):1787–1796, October 2023.

[128] Samuel Westby and Christoph Riedl. Collective Intelligence in Human-AI Teams: A Bayesian Theory of Mind Approach. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(5):6119–6127, June 2023.

[129] Rhyse Bendell, Jessica Williams, Stephen M. Fiore, and Florian Jentsch. Individual and team profiling to support theory of mind in artificial social intelligence. *Scientific Reports*, 14(1):12635, June 2024.

[130] Ştefan Sarkadi. An Arms Race in Theory-of-Mind: Deception Drives the Emergence of Higher-level Theory-of-Mind in Agent Societies. In *2023 IEEE International Conference on Autonomic Computing and Self-Organizing Systems (ACSOS)*, pages 1–10, September 2023.

[131] Te Bao, Sascha Füllbrunn, Jiaoying Pei, and Jichuan Zong. Reading the market? Expectation coordination and theory of mind. *Journal of Economic Behavior & Organization*, 219:510–527, March 2024.

[132] Alvaro Garcia-Lopez. Theory of Mind in Artificial Intelligence Applications. In Teresa Lopez-Soto, Alvaro Garcia-Lopez, and Francisco J. Salguero-Lamillar, editors, *The Theory of Mind Under Scrutiny: Psychopathology, Neuroscience, Philosophy of Mind and Artificial Intelligence*, pages 723–750. Springer Nature Switzerland, Cham, 2023.

[133] Michele Rocha, Heitor Henrique da Silva, Analúcia Schiaffino Morales, Stefan Sarkadi, and Alison R. Panisson. Applying Theory of Mind to Multi-agent Systems: A Systematic Review. In Murilo C. Naldi and Reinaldo A. C. Bianchi, editors, *Intelligent Systems*, pages 367–381, Cham, 2023. Springer Nature Switzerland.

[134] Alberto Nebreda, Danylyna Shpakivska-Bilan, Carmen Camara, and Gianluca Susi. The Social Machine: Artificial Intelligence (AI) Approaches to Theory of Mind. In Teresa Lopez-Soto, Alvaro Garcia-Lopez, and Francisco J. Salguero-Lamillar, editors, *The Theory of Mind Under Scrutiny: Psychopathology, Neuroscience, Philosophy of Mind and Artificial Intelligence*, pages 681–722. Springer Nature Switzerland, Cham, 2023.

[135] Edmond Awad, Sohan Dsouza, Richard Kim, Jonathan Schulz, Joseph Henrich, Azim Shariff, Jean-François Bonnefon, and Iyad Rahwan. The moral machine experiment. *Nature*, 563(7729):59–64, 2018.

[136] Yuanyuan Mao, Shuang Liu, Qin Ni, Xin Lin, and Liang He. A review on machine theory of mind. *IEEE Transactions on Computational Social Systems*, 2024.

[137] Javier Fernández-López de Pablo, Valéria Romano, Maxime Derex, Erik Gjesfjeld, Claudine Gravel-Miguel, Marcus J Hamilton, Andrea Bamberg Migliano, Felix Riede, and Sergi Lozano. Understanding hunter–gatherer cultural evolution needs network thinking. *Trends in Ecology & Evolution*, 37(8):632–636, 2022.

[138] Mark Dyble, James Thompson, Daniel Smith, Gul Deniz Salali, Nikhil Chaudhary, Abigail E Page, Lucio Vinicuis, Ruth Mace, and Andrea Bamberg Migliano. Networks of food sharing reveal the functional significance of multilevel sociality in two hunter-gatherer groups. *Current Biology*, 26(15):2017–2021, 2016.

[139] Kevin Laland, Blake Matthews, and Marcus W Feldman. An introduction to niche construction theory. *Evolutionary ecology*, 30:191–202, 2016.

[140] Peter Rowley-Conwy and Robert Layton. Foraging and farming as niche construction: stable and unstable adaptations. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 366(1566):849–862, 2011.

[141] Manuel Arroyo-Kalin, David Wengrow, Dorian Q Fuller, Chris J Stevens, and Michèle Wollstonecroft. Civilisation and human niche construction. *Archaeology International*, 20(1):106–109, 2017.

[142] Melissa E Kemp, Alexis M Mychajliw, Jenna Wadman, and Amy Goldberg. 7000 years of turnover: historical contingency and human niche construction shape the caribbean's anthropocene biota. *Proceedings of the Royal Society B*, 287(1927):20200447, 2020.

[143] Bruce D Smith and Melinda A Zeder. The onset of the anthropocene. *Anthropocene*, 4:8–13, 2013.

[144] Erle C Ellis. The anthropocene condition: evolving through social–ecological transformations. *Philosophical Transactions of the Royal Society B*, 379(1893):20220255, 2024.

[145] David Sloan Wilson, Guru Madhavan, Michele J Gelfand, Steven C Hayes, Paul WB Atkins, and Rita R Colwell. Multilevel cultural evolution: From new theory to practical applications. *Proceedings of the National Academy of Sciences*, 120(16):e2218222120, 2023.

[146] Paul A Ryan, Simon T Powers, and Richard A Watson. Social niche construction and evolutionary transitions in individuality. *Biology & philosophy*, 31:59–79, 2016.

[147] A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.

[148] David Silver, Satinder Singh, Doina Precup, and Richard S Sutton. Reward is enough. *Artificial Intelligence*, 299:103535, 2021.

[149] Zhiyuan Li, Hong Liu, Denny Zhou, and Tengyu Ma. Chain of thought empowers transformers to solve inherently serial problems. *arXiv preprint arXiv:2402.12875*, 2024.

[150] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.

[151] James WA Strachan, Dalila Albergo, Giulia Borghini, Oriana Pansardi, Eugenio Scaliti, Saurabh Gupta, Krati Saxena, Alessandro Rufo, Stefano Panzeri, Guido Manzi, et al. Testing theory of mind in large language models and humans. *Nature Human Behaviour*, pages 1–11, 2024.

[152] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.

[153] Momiao Xiong. *Artificial Intelligence and Causal Inference*. Chapman and Hall/CRC, 2022.

[154] Kevin Laland, Tobias Uller, Marc Feldman, Kim Sterelny, Gerd B Müller, Armin Moczek, Eva Jablonka, John Odling-Smee, Gregory A Wray, Hopi E Hoekstra, et al. Does evolutionary theory need a rethink? *Nature*, 514(7521):161–164, 2014.

[155] Eörs Szathmáry and John Maynard Smith. The major evolutionary transitions. *Nature*, 374(6519):227–232, 1995.

[156] Eörs Szathmáry. Toward major evolutionary transitions theory 2.0. *Proceedings of the National Academy of Sciences*, 112(33):10104–10111, 2015.

[157] Mikhail Prokopenko, Paul CW Davies, Michael Harré, Marcus Heisler, Zdenka Kuncic, Geraint F Lewis, Ori Livson, Joseph T Lizier, and Fernando E Rosas. Biological arrow of time: Emergence of tangled information hierarchies and self-modelling dynamics. *arXiv preprint arXiv:2409.12029*, 2024.

[158] David Sloan Wilson and Dennis J Snower. Rethinking the theoretical foundation of economics i: The multilevel paradigm. *Economics*, 18(1):20220070, 2024.

[159] Dániel Czégel, István Zachar, and Eörs Szathmáry. Multilevel selection as bayesian inference, major transitions in individuality as structure learning. *Royal Society open science*, 6(8):190202, 2019.

# A  Utility polynomials

In two-player, $m$-action normal-form games, the payoffs are encoded in a single $m \times m$ matrix $\mathtt{G}$, where each cell contains the payoffs for the $i$ and $j$ players for a particular combination of actions.

|       |     | $x_j$ | |
|-------|-----|-------|-------|
|       |     | C     | D     |
| $x_i$ | C   | $(g_i^{cc}, g_j^{cc})$ | $(g_i^{cd}, g_j^{dc})$ |
|       | D   | $(g_i^{dc}, g_j^{cd})$ | $(g_i^{dd}, g_j^{dd})$ |

Table 2: Two-player, two-action normal-form game payoff matrix $\mathtt{G}$

The payoff matrix $\mathtt{G}$ can be decomposed into two standard matrices $\mathtt{G}_i, \mathtt{G}_j$ containing each player's respective payoffs. Utility functions can be obtained in polynomial form from the elements $g_i^l \in \mathtt{G}_i$, where each term intuitively represents the contribution of each agent's action (first-order terms) and their interaction (second-order terms) [117]. For $m = 2$ actions, we have:

$$U_i(\mathbf{x}; \mathbf{a}_i) = a_i^0 + a_i^i x_i + a_i^j x_j + a_i^{ij} x_i x_j \tag{A1}$$

$$\mathbf{a}_i = \begin{bmatrix} a_i^0 \\ a_i^i \\ a_i^j \\ a_i^{ij} \end{bmatrix} = \frac{1}{4} \begin{bmatrix} g_i^{cc} + g_i^{cd} + g_i^{dc} + g_i^{dd} \\ (g_i^{dc} + g_i^{dd}) - (g_i^{cc} + g_i^{cd}) \\ (g_i^{cd} + g_i^{dd}) - (g_i^{cc} + g_i^{dc}) \\ (g_i^{cc} + g_i^{dd}) - (g_i^{cd} + g_i^{dc}) \end{bmatrix} \tag{A2}$$

## A.1  General form for $n$ players and $m = 2$ actions

For any number $n$ of players where each player's action $x_i \in \{-1, 1\}$, the utility function is

$$U_i(\mathbf{x}; \mathbf{a}_i) = \sum_{S \subseteq \{1,2,\ldots,n\}} a_i^S \prod_{j \in S} x_j \tag{A3}$$

where $S$ is any subset of the set $\{1, 2, \ldots, n\}$, and $a_i^S$ is the co-factor corresponding to the subset $S$. For example, $S = \{1, 2, 3\}$ means the product $\prod_{j \in S} x_j = x_1 x_2 x_3$, and $a_i^S = a_i^{1,2,3}$ represents the co-factor for this interaction. To determine the co-factors $a_i^S$, we need to consider the payoff matrix $\mathtt{G}_i$, which contains $2^n$ entries corresponding to each combination of $n$ players' actions. Denote the entry in $\mathtt{G}_i$ for the action combination $(x_1, x_2, \ldots, x_n)$ as $g_i^{x_1 x_2 \ldots x_n}$. The co-factors $a_i^S$ are:

$$a_i^S = \frac{1}{2^n} \sum_{(x_1, x_2, \ldots, x_n) \in \{-1, 1\}^n} g_i^{x_1 x_2 \ldots x_n} \prod_{j \in S} x_j \tag{A4}$$

This generalization for the utility function for $n$-player, two-action games can be interpreted through Fourier analysis on the Boolean cube. This connection arises because we represent the utility function as a sum of basis functions over the Boolean domain, where the co-factors $a_i^S$ are Fourier co-factors and $\prod_{j \in S} x_j$ are parity functions. Thus, we express the utility function $U_i$ as a Fourier transform of the payoff matrix $\mathtt{G}_i$. Each co-factor $a_i^S$ captures the influence of a subset $S$ of players on the overall utility, analogous to how Fourier co-factors capture the influence of different frequency components in a signal. This generalization leverages the principles of Fourier analysis on the Boolean cube to decompose the complex payoff interactions into simpler, orthogonal components, making it easier to analyze and interpret the contributions of different action combinations in the game.

# B    Derivation of utility polynomials in example

From the provided payoff matrix,

|   | C | D |
|---|---|---|
| C | $(1,1)$ | $(0-c, 1+c)$ |
| D | $(1+c, 0-c)$ | $(0,0)$ |

we can obtain the polynomial co-factors (see Table 2 and Equation A2 for details)

$$\mathbf{a}_i = \begin{bmatrix} a_i^0 \\ a_i^i \\ a_i^j \\ a_i^{ij} \end{bmatrix} = \frac{1}{4} \begin{bmatrix} g_i^{cc} + g_i^{cd} + g_i^{dc} + g_i^{dd} \\ (g_i^{dc} + g_i^{dd}) - (g_i^{cc} + g_i^{cd}) \\ (g_i^{cd} + g_i^{dd}) - (g_i^{cc} + g_i^{dc}) \\ (g_i^{cc} + g_i^{dd}) - (g_i^{cd} + g_i^{dc}) \end{bmatrix} \tag{B1}$$

$$= \frac{1}{4} \begin{bmatrix} 1 + (0-c) + (1+c) + 0 \\ ((1+c)+0) - (1+(0-c)) \\ ((0-c)+0) - (1+(1+c)) \\ (1+0) - ((0-c)+(1+c)) \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 1+1 \\ 1+c-1+c \\ -c-1-1-c \\ 1+c-1-c \end{bmatrix} \tag{B2}$$

$$= \frac{1}{4} \begin{bmatrix} 2 \\ 2c \\ -2(1+c) \\ 0 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 1 \\ c \\ -(1+c) \\ 0 \end{bmatrix} = \begin{bmatrix} a_i^0 \\ a_i^i \\ a_i^j \\ a_i^{ij} \end{bmatrix} \tag{B3}$$

and substitute into the (second-degree) utility polynomial (with zero quadratic term in this particular case)

$$U_i(\mathbf{x}; \mathbf{a}) = a_i^0 + a_i^i x_i + a_i^j x_j + a_i^{ij} x_i x_j = \frac{1}{2} \Big( 1 + c x_i - (1+c) x_j \Big) \tag{B4}$$

we can choose to re-scale the utility values by a factor of 2, which results in the following utility polynomials for $A_2$ and $A_3$, respectively

$$U_2(\mathbf{x}; \mathbf{a}) = 1 + c x_2 - (1+c) x_3, \qquad\qquad U_3(\mathbf{x}; \mathbf{a}) = 1 + c x_3 - (1+c) x_2. \tag{B5}$$