



Université Libre de Bruxelles

*Institut de Recherches Interdisciplinaires
et de Développements en Intelligence Artificielle*

**Learning from humans to build social
cognition among robots**

N. COUCKE, M.K. HEINRICH, A. CLEEREMANS, and M.
DORIGO

IRIDIA – Technical Report Series

Technical Report No.
TR/IRIDIA/2022-008

September 2022
Last revision: January 2023

IRIDIA – Technical Report Series
ISSN 1781-3794

Published by:

IRIDIA, *Institut de Recherches Interdisciplinaires
et de Développements en Intelligence Artificielle*
UNIVERSITÉ LIBRE DE BRUXELLES
Av F. D. Roosevelt 50, CP 194/6
1050 Bruxelles, Belgium

Technical report number TR/IRIDIA/2022-008

Revision history:

TR/IRIDIA/2022-008.001	September 2022
TR/IRIDIA/2022-008.002	January 2023

The information provided is the sole responsibility of the authors and does not necessarily reflect the opinion of the members of IRIDIA. The authors take full responsibility for any copyright breaches that may result from publication of this paper in the IRIDIA – Technical Report Series. IRIDIA is not responsible for any use that might be made of data appearing in this publication.

Learning from humans to build social cognition among robots

Nicolas Coucke^{1,2,†,*}, Mary Katherine Heinrich^{1,†,*}, Axel Cleeremans²
and Marco Dorigo¹

¹*IRIDIA, Université Libre de Bruxelles, Brussels, Belgium*

²*Consciousness, Cognition & Computation Group, Center for Research in Cognition and Neurosciences, Université Libre de Bruxelles, Brussels, Belgium*

[†] *These authors contributed equally to this work and share first authorship.*

Correspondence*:

Nicolas Coucke
nicolas.coucke@ulb.be

Mary Katherine Heinrich
mary.katherine.heinrich@ulb.be

2 ABSTRACT

3 Self-organized groups of robots have generally coordinated their behaviors using quite simple
4 social interactions. Although simple interactions are sufficient for some group behaviors, future
5 research needs to investigate more elaborate forms of coordination, such as social cognition, to
6 progress towards real deployments. In this perspective, we define social cognition among robots
7 as the combination of social inference, social learning, social influence, and knowledge transfer,
8 and propose that these abilities can be established in robots by building underlying mechanisms
9 based on behaviors observed in humans. We review key social processes observed in humans
10 that could inspire valuable capabilities in robots and propose that relevant insights from human
11 social cognition can be obtained by studying human-controlled avatars in virtual environments
12 that have the correct balance of embodiment and constraints. Such environments need to allow
13 participants to engage in embodied social behaviors, for instance through situatedness and
14 bodily involvement, but, at the same time, need to artificially constrain humans to the operational
15 conditions of robots, for instance in terms of perception and communication. We illustrate our
16 proposed experimental method with example setups in a multi-user virtual environment.

17 **Keywords:** artificial social cognition, embodied cognition, self-organization, robot swarms, multi-robot systems, artificial intelligence,
18 artificial general intelligence, social robots

1 INTRODUCTION

19 AI research has greatly advanced, but when interaction with other agents is required, existing algorithms
20 easily break down (Bard et al., 2020). Social interaction and social embodiment are still underexplored in
21 artificial general intelligence (Bolotta and Dumas, 2022) and in groups of intelligent robots. While there is
22 some robotics research on social cognition, it focuses on human-robot interaction (Henschel et al., 2020),
23 e.g., how a robot interprets the intentions of a human, not on interactions *among* robots. It is important
24 to note that what looks like social cognition is not necessarily social cognition. For instance, agents or
25 robot controllers made by reinforcement learning might behave in ways that look socially cognizant in

26 some situations, but this might only be appearance—i.e., the underlying behavioral phenomena are not
27 there—so the illusion will break down when exposed to more situations.

28 Robots can coordinate with each other by using, e.g., centralized control or self-organization. In multi-
29 robot systems that are not self-organized, robots are directed to follow a centrally coordinated plan using
30 explicit commands or global references. In this paper, we are interested exclusively in robot groups
31 that include aspects of self-organization, because social cognition depends on some degree of individual
32 autonomy. If a robot is essentially a remote-controlled sensor or actuator, it does not engage in social
33 cognition.

34 In existing research on self-organized robot groups, the individuals are usually quite simple and often
35 rely on indiscriminate, naïve interactions. Indeed, swarm robotics research has shown that no advanced
36 cognition or elaborate social negotiation is needed to self-organize certain group behaviors (e.g., Nouyan
37 et al., 2009; Rubenstein et al., 2014; Valentini et al., 2016). However, it has been argued that there are still
38 significant gaps for robot swarms to be deployment-ready, and that the future of swarm robotics research
39 should concentrate on more elaborate forms of self-organized coordination (Dorigo et al., 2020, 2021), such
40 as self-organized hierarchy (Mathews et al., 2017; Zhu et al., 2020) or behavioral heterogeneity (Kengyel
41 et al., 2015).

42 In this perspective, we argue that another important direction for future study should be social cognition.
43 Robot groups successfully equipped with social cognition could engage in elaborate coordination without
44 sending each other large amounts of data. Some aspects of robot behavior could be mutually predictable,
45 for instance by robots maintaining good internal models of each other. Socially cognitive robots could have
46 improved group performance, e.g., by not destructively interfering with each other (which requires time
47 and effort to resolve) and not accidentally disrupting each other’s sub-goals while attempting to reach a
48 common goal.

49 In cognitive robotics, research on individual robots such as humanoids is very advanced (Cangelosi
50 and Asada, 2022), even on each of the six key attributes of artificial cognitive systems (Vernon, 2014):
51 action, perception, autonomy, adaptation, learning, and anticipation. Comparatively, cognition in swarm
52 robotics is still in its beginning stages. While cognitive robot swarms can be autonomously capable of
53 collective action, perception, and in some cases adaptation (Heinrich et al., 2022), we do not yet know how
54 to make robot swarms that can autonomously learn and anticipate as a collective, in such a way that the
55 group behavior is greater than the sum of its parts. We propose that studying social cognition could help us
56 advance the autonomous collective capabilities of groups of robots.

2 SOCIALLY COGNITIVE ROBOTS: OUR PERSPECTIVE

57 Our perspective is summarized as follows: social cognition among robots can be built by developing
58 artificial social reasoning capabilities based on behaviors observed in humans.

59 Frith (2008) has defined social cognition in humans as “the various psychological processes that enable
60 individuals to take advantage of being part of a social group” and Frith and Frith (2012) have further
61 specified that a substantial portion of these psychological processes are for learning about and making
62 predictions about other members of the social group. The mechanisms of social cognition in humans
63 include social signalling, social referencing, mentalizing (i.e., tracking of others’ mental states, intended
64 actions, objectives, and opinions), observational learning (e.g., social reward learning, mirroring), deliberate
65 knowledge transfer (e.g., teaching), and sharing of experiences through reflective discussion (Frith, 2008;

66 Frith and Frith, 2012). Crucially, social cognition is also defined as “not reducible to the workings of
67 individual cognitive mechanisms” (De Jaegher et al., 2010).

68 Although some social abilities such as simple social interaction are well-developed among robots, most
69 of the abilities contained in Frith (2008)’s definition of social cognition are lacking, and could provide
70 significant performance benefits. For instance, the transfer of information between robots is well understood,
71 but much less so the transfer of knowledge, especially implicitly.

72 We define social cognition among robots as the following set of abilities:

- 73 1. **social inference** – inferring the opinions, intended next actions, and overall goals of other robots in the
74 same social group, using interpretation of social signals;
- 75 2. **social learning** – learning information about which actions to adopt or avoid based on observations of
76 each other’s behaviors and social signalling;
- 77 3. **social influence** – deliberately influencing each other’s (socially inferred) internal states using social
78 signaling; and
- 79 4. **knowledge transfer** – transferring high-level knowledge using social interaction, e.g., using implicit
80 demonstration or explicit instruction.

81 Currently, robots are well-equipped with some of the requirements for these abilities, such as simple social
82 interactions, but lack other crucial requirements such as explicit social reasoning. Although research has
83 shown that no social cognition is needed for simple group behaviors in robots, it is an open challenge how
84 to accomplish more advanced behaviors in a fully self-organized way. Some of the significant unresolved
85 technical challenges for advanced self-organization among robots, which we believe social cognitive
86 abilities could contribute to, are the following:

- 87 • autonomously anticipating which actions should be taken in an environment filled with other
88 autonomous robots,
- 89 • collectively defining an explicit goal that was not pre-programmed and collectively directing the robot
90 group towards it,
- 91 • making online inferences about other robots’ current states and future behaviors, and adapting
92 their coordination strategies accordingly, even while moving at high speed in dynamic unknown
93 environments, and
- 94 • designing self-organization among robots such that the resulting group behaviors, although not
95 completely predictable, are safe and trustable.

96 We propose that socially cognitive robots can in part be developed by learning from the social cognition
97 processes of humans in certain experimental conditions. In order to have the potential to transfer observed
98 behaviors and capabilities from humans to robots, we believe experiments with human subjects must be
99 conducted in a platform that allows experimental setups to be: on one hand, realistic enough to study
100 **embodied** human behavior, but on the other hand, **constrained** and simplified enough to approximate the
101 operational conditions of robots.

3 STATE OF THE ART

102 3.1 Artificial social learning and artificial mentalizing

103 Many examples of artificial learning exist that seem relevant to the mechanisms of social cognition.
104 However, key social aspects are not present in these existing methods: for instance, reward learning
105 has been demonstrated in robots (e.g., Daniel et al., 2015) but learning of social rewards among robots
106 has not been studied. Likewise, robots learning by interacting with and observing other robots has been
107 demonstrated (e.g., Murata et al., 2015), but not for the learning of socially relevant information nor to
108 build behaviors among robots that are irreducible to the knowledge held by robots individually.

109 Currently, the most advanced research towards artificial social cognition can be seen in multi-agent
110 reinforcement learning. In basic approaches, each agent would use reinforcement learning individually,
111 treating other agents as part of the environment. In more elaborate existing approaches, agents are trained
112 to model each other and several types of artificial mentalizing have been demonstrated (Albrecht and
113 Stone, 2018). For example, in the Deep Reinforcement Opponent Network (DRON), one agent learns
114 the representation of the opponent’s policy (He et al., 2016). In another example, an agent uses itself
115 as the basis to predict another agent’s actions (Raileanu et al., 2018). One approach using a “Theory of
116 Mind” network has even produced agents that can explicitly report inferred mental states of other agents
117 and passes the classic “false belief test” for understanding the mental states of others (Rabinowitz et al.,
118 2018). Current efforts in multi-agent learning use cooperative games such as Hanabi as benchmarks, which
119 involves inferring other’s mental state and using that information to collaborate (Bard et al., 2020). For
120 the development of artificial social cognition, the next step for this line of research would be to situate the
121 mentalizing behaviors within the full set of social cognition mechanisms, including social influence and
122 social reward learning (cf. Olsson et al., 2020).

123 3.2 Social cognition transfer between humans and robots

124 Robots have been used as experimental tools for the study of embodied social cognition. For instance,
125 a variety of devices have been used to automatically provide synthetic social stimuli to animals in a
126 naturalistic way (Frohnwieser et al., 2016). Similarly, the effect of humanoid robots on human social
127 cognition has been broadly studied (Wykowska et al., 2016). Social robots in the context of human-robot
128 interaction have also been investigated (e.g., Dautenhahn, 2007). However, to the best of our knowledge,
129 no studies have looked at expanding these robot use cases into embodied artificial social cognition among
130 robots, and no work apart from our own has proposed using experiments with humans to contribute to
131 building social cognition among robots.

4 DIRECTIONS FOR FUTURE RESEARCH

132 Advanced group capabilities seen in humans can inspire similar capabilities in robots. For example, the
133 human capabilities of selecting and following leaders (Van Vugt, 2006) and re-organizing communication
134 networks around individuals with better information (Almaatouq et al., 2020) have recently inspired the
135 development of self-organized hierarchies for robots, for instance using physical (Mathews et al., 2017) or
136 wireless connections (Zhu et al., 2020). In the following sections, we identify cognitive processes used by
137 humans in social situations that would be valuable for robot groups, and propose them as future research
138 directions for building social cognition among robots.

139 4.1 Social heuristics and action selection

140 Humans often use cognitive processes known as “heuristics” to select actions in social situations. In
141 humans, heuristics are defined as action selection strategies that usually deviate from economic rationality
142 or Bayesian optimality but which facilitate a rapid action selection when time and knowledge about a
143 situation are limited (Hertwig and Herzog, 2009). The hidden states of other agents cannot be directly
144 observed, so the outcome of a social situation always has a high degree of uncertainty—selecting the
145 optimal action is computationally intractable (Seymour and Dolan, 2008).

146 In humans, heuristics can involve continuous integration of multiple variables or sources of
147 information, for example when deciding on a walking direction based on the position of other walking
148 individuals (Moussaid et al., 2011). In psychology and neuroscience, action selection is often characterized
149 as the result of an accumulation process, in which evidence that supports a certain decision or action is
150 accumulated over time (Ratcliff and McKoon, 2008). A certain action is taken when the accumulated
151 evidence crosses some threshold. The sources and manner of evidence integration can be determined by
152 social heuristics. For example, evidence accumulation frameworks can characterize how humans use a
153 “follow the majority heuristic” during social decision making (Tump et al., 2020), as well as how humans
154 base their own movements on those of others during embodied competitive interactions (Lokesh et al.,
155 2022).

156 4.2 Coupling, alignment, and mirroring

157 Humans often mirror each other’s behaviors and can participate in a “coupling” behavior through
158 reciprocal interactions. Implicit coupling can occur between physiological states (for example,
159 synchronization of heartbeats and breathing rhythms). Explicit sensorimotor coupling involves mutual
160 prediction of each other’s actions and facilitates coordinated action sequences (Dumas and Fairhurst, 2021).
161 On a higher cognitive level, reciprocal interactions can create alignment between internal cognitive states,
162 which in turn facilitates better mutual prediction of actions (Friston and Frith, 2015).

163 Humans can also disengage from social interactions and instead mirror (or “simulate”) others’ actions as a
164 type of internalized action (Buzsáki, 2019, p. 131). This capacity is supported by the mirror neuron system,
165 which is active when observing and when executing a movement (Rizzolatti and Craighero, 2004). Internal
166 simulation aids in understanding others’ intentions and in selecting complementary actions (Newman-
167 Norlund et al., 2007).

168 4.3 Mentalizing and shared representations

169 Simply mirroring the mental states of others is often not sufficient to infer their opinions, objectives, or
170 intended actions (Saxe, 2005). Therefore, coupling and mirroring are often complemented in humans by
171 higher-level cognition about others’ beliefs, desires, and intentions, taking into account factors such as
172 context and memory (Sebanz et al., 2006). This requires mentalizing, a process of inference about others’
173 changing mental states, beyond simple mirroring (Frith and Frith, 2012).

174 For example, mentalizing based on observations of others’ gazes facilitates taking others’ perspectives
175 into account and tracking their beliefs about a shared environment or world (Frith and Frith, 2012). By
176 observing others’ movements, humans can also infer the confidence that others have in their beliefs (Patel
177 et al., 2012) and the intentions that underlie their actions (Baker et al., 2009). Crucially, humans also
178 mentalize based on third-party observations of others’ interactions, and then estimate the social relationships
179 between them (Ullman et al., 2009).

180 Tracking others' goals and beliefs helps humans to distinguish which subset of their action representations
181 are shared with others. Shared representations aid in predicting and interpreting the actions of others in
182 the context of a joint goal, and in selecting complementary actions. For instance, by tracking others'
183 beliefs, an individual can recognize when communication or signalling is needed to facilitate smooth
184 coordination (Pezzulo and Dindo, 2011).

185 **4.4 Outcome monitoring**

186 Humans monitor behaviors and detect errors when taking actions directed towards a certain
187 goal (Botvinick et al., 2001). If an individual recognizes another making what might be an error, in
188 pursuit of a shared goal, the individual needs to then distinguish whether it was indeed an error, or whether
189 their goals are misaligned.

190 Humans also monitor whether actions have their intended outcomes, as well as whether a certain action
191 and certain outcome actually have a causal link. This results in a greater or lesser sense of agency over a
192 certain action or outcome (Haggard and Chambon, 2012), which in turn impacts how an individual acts
193 in social situations. Agency can be modulated in a variety of ways: joint agency when acting together
194 with others, vicarious agency when influencing the actions of others, or violated agency when actions are
195 interfered with by others (Silver et al., 2020). The modulated sense of agency in humans helps shape an
196 individual's monitoring of links between actions, errors, and outcomes.

5 FROM HUMANS TO ROBOTS: AN EXPERIMENTAL METHOD

197 Robots are embodied agents with specific morphologies and specific perception and action capabilities
198 that differ from (and are often far more limited than) those of humans. To gain insights from human social
199 cognition that are relevant to robots, human subjects would need to be studied in an experimental platform
200 that: (i) allows them to engage in embodied social behaviors, but also (ii) allows enough constraints to
201 artificially expose humans to the operational conditions of robots. We propose that behavioral experiments
202 conducted with humans controlling avatars in virtual environments can achieve this trade-off.

203 **5.1 Balancing embodiment and constraints in virtual environments**

204 Existing experiments on human social cognition have mostly been conducted in highly controlled single-
205 person paradigms which lack embodiment. We identify the following five aspects of embodiment that we
206 propose human-controlled avatars in new virtual environments will need, for the study of embodied human
207 social cognition.

- 208 1. **Situatedness:** An agent takes actions while being part of a situation, rather than by observing the
209 situation from the outside (Wilson, 2002).
- 210 2. **Sensory and action shaping:** By taking actions (e.g., moving their bodies) in the environment,
211 agents can actively change the flow of their sensory inputs as well as the potential effects of their
212 actions (Gordon et al., 2021).
- 213 3. **Bodily involvement:** The bodily state and/or morphology of the agent—as well as the agent's bodily
214 relation to the bodies of other agents—can be involved in cognition (Wilson, 2002).
- 215 4. **Interaction cascades:** Agents can engage with each other in such a way that actions by one can
216 influence reciprocal actions by another, resulting in cascades of interactions and behaviors (Dale et al.,
217 2013).

218 5. **High bandwidth:** There can be high bandwidth of implicit or explicit information exchange between
219 agents (Schilbach et al., 2013).

220 Complementarily, we identify the following constraints that will also need to be possible in the virtual
221 environment.

222 1. **Body and action:** Human-controlled avatars can be equipped with morphology features and action
223 capabilities that are similar to those of relevant robots.

224 2. **Perception:** When controlling an avatar, a human subject can be limited to sensory inputs similar to
225 those of relevant robots (e.g., restricted visual information).

226 3. **Communication:** Human-controlled avatars can be limited to communication and signalling
227 capabilities that are similar to those available to relevant robots.

228 4. **Hidden states:** Human subjects can be required to explicitly report information about hidden states
229 (e.g., their current opinion or confidence level) that is not directly observable from their behavior but
230 would be available to an experimenter if using relevant robots.

231 Unconstrained real-world social situations would fulfill all listed requirements for embodiment, but would
232 lack control and interpretability. Virtual environments enable certain aspects of embodiment while at the
233 same time ensuring a degree of control of the situation for the experimenter.

234 5.2 Example: using the virtual environment HuGoS

235 To the best of our knowledge, no off-the-shelf virtual environment was available to meet these
236 requirements, so we built a tool in Unity3D called “HuGoS: Humans Go Swarming” (Coucke et al.,
237 2020, 2021) that we could use to study human behavior in embodied scenarios similar to those in which
238 robots operate. To illustrate the features that we propose for a virtual environment for studying transferable
239 social cognition, we describe two example experimental setups in HuGoS.

240 5.2.1 Collective decision making

241 Collective decision making has been widely studied in swarm robotics (Valentini et al., 2017), but
242 many gaps still remain (Khaluf et al., 2019). Collective decisions have also been extensively studied in
243 humans (Kameda et al., 2022), but not typically in embodied scenarios that would be relevant to robots,
244 in which, e.g., exploration and signalling can take place simultaneously. In our example implementation
245 in Coucke et al. (2020), each of four participants controls the movements of a cubic avatar in an environment
246 scattered with red and blue cylindrical landmarks (see Fig. 1). The task is to explore the environment while
247 making observations through the avatar’s (broad or limited) field of view and simultaneously deciding
248 whether there are more red or blue landmarks present in the environment. The participants must come to a
249 consensus in order to complete the task and are only permitted to communicate with each other indirectly:
250 they vote by changing their avatar color and they observe the avatar colors of the other participants while
251 making their decisions (see Fig. 1a-c). During an experiment, all perceptual information available to each
252 participant, along with their actions, are recorded at a sampling rate of 10 Hz (Fig. 1d-h).

253 In this experiment setup, participants came to a consensus about the predominant color in the environment
254 through a combination of environmental and social information. In the example trial shown in Fig. 1, at 45 s,
255 all four participants had adopted the correct opinion (Fig. 1h) after individually and broadly exploring the
256 environment and then reducing their average relative distances to increase their access to social information
257 (Fig. 1f) and finally come to a consensus. When a consensus was reached, not all participants had personally
258 observed all parts of the environment (Fig. 1e), implying that social information was effectively used.

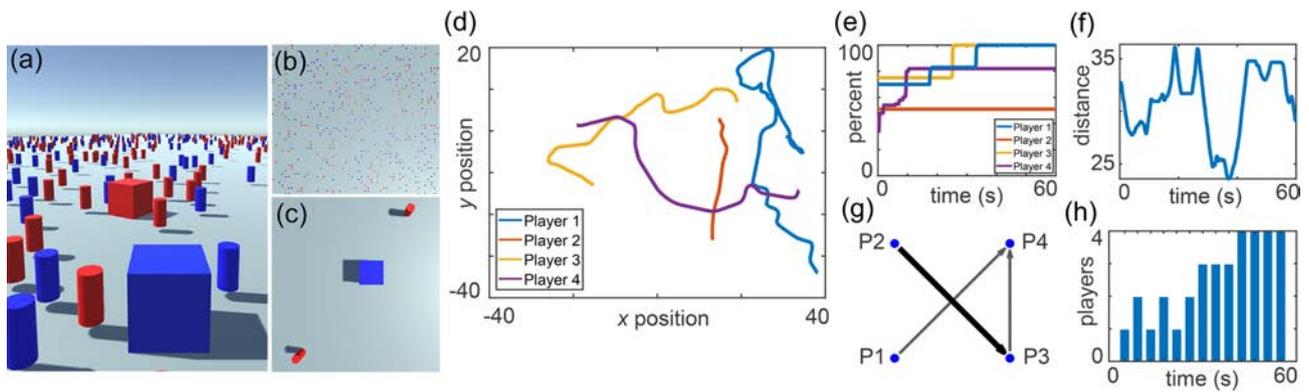


Figure 1. Collective decision making. Participants control cubic avatars while having either a broad (a) or limited (c) view of the full environment (b). A wide variety of variables can be measured during the experiment, such as the participants' trajectories (d), the percentage of the environment they have explored (e), the average distance between participants (f), the participant-participant viewing network (g), and the number of avatars choosing the correct color (h). Figure reprinted by permission from Springer Nature Customer Service Centre GmbH: Springer eBook, Coucke et al. (2020), © Springer Nature 2020.

259 Further, all participants had a strong directional line-of-sight connection with at least one other participant
 260 (Fig. 1g), but the most looked-at participant (P4) had not personally observed the whole environment
 261 (Fig. 1e), implying that the consensus on the correct opinion was indeed arrived at by a self-organized
 262 and collective process. For more information on this and similar experiments, please refer to Coucke et al.
 263 (2020). By setting up more advanced experiments in this direction, data could be collected to, for example,
 264 identify social heuristics that can inspire new protocols in future robot swarms.

265 5.2.2 Collective construction

266 Existing swarm robotics approaches to construction often use stigmergy (i.e., indirect communication
 267 through modification of the environment) to coordinate (Petersen et al., 2019), but the structures built
 268 strictly by stigmergy are relatively simple. Future robot swarms should be able to build complex structures
 269 in dynamically changing environments (Dorigo et al., 2020). In our example 'lava spill task' scenario
 270 in Coucke et al. (2021), human social behaviors in collective construction scenarios can be observed. In this
 271 task (see Fig. 2), participants are instructed to collectively construct a barrier to contain an expanding spill,
 272 but are not instructed how to coordinate. Each participant controls the movement of an avatar that can push
 273 construction blocks. The environment includes two different spills (i.e., expanding circles) and a supply of
 274 construction blocks placed in between them. During an experiment, a group of eight participants needs
 275 to assess the environment and coordinate their actions using indirect communication (i.e., observation of
 276 peers) to barricade both of the expanding spills within 300 seconds.

277 The avatar trajectories in Fig. 2e show that participants coordinated to distribute their work between the
 278 two spills and place blocks around the full circumferences of both spills. Fig. 2d shows that participants
 279 continued to place more blocks at a roughly constant rate throughout the experiment, implying that no
 280 bottleneck arose in their self-organized coordination. The figure also shows that the expansion of both
 281 spills had successfully been stopped at around 200 s. For more information on this and similar experiments,
 282 please refer to Coucke et al. (2021). Using more advanced setups in this direction, the gathered behavioral
 283 data could provide insights into how self-organized coordination and group actions unfold over time and
 284 adapt to the environment. In order to get detailed information about participant strategies, experiments in
 285 this virtual environment can be temporarily interrupted at certain times to ask participants about, e.g., their

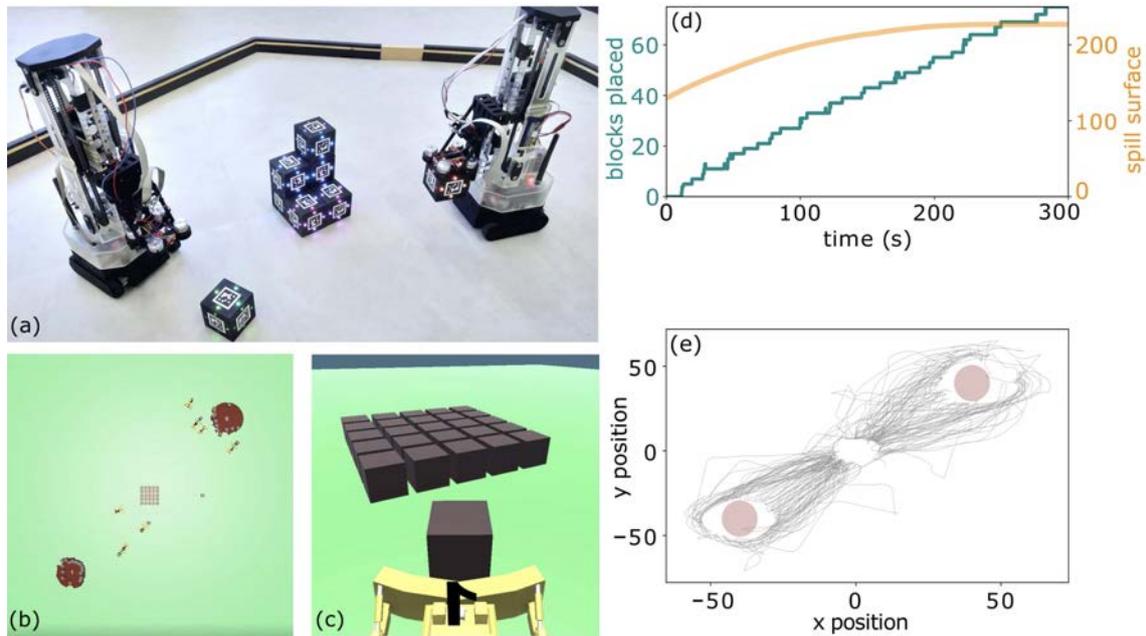


Figure 2. Collective construction. **(a)** Two physical robots that perform collective construction using stigmergic blocks (Allwright et al., 2019). Figure a reprinted from Allwright et al. (2019) under license CC BY-NC-ND 4.0. **(b-c)** ‘Lava spill task’ in which participants use indirect communication to collectively construct a barrier to contain expanding spills. **(d)** The spill size stagnates after around 200 seconds, when participants successfully enclosed it with construction blocks. **(e)** Data such as the avatar trajectories can be used to analyze how participants coordinate the placement of blocks. Figures b-d adapted from Coucke et al. (2021) under license CC BY 4.0.

286 explicit judgements about the beliefs of other participants, their sense of (joint) agency, or their feelings of
 287 alignment with others.

6 DISCUSSION

288 Some features of human social groups, such as collective intentions, reflective discussion, or shared
 289 biases, might at first seem not particularly relevant for robots. However, there are many autonomous
 290 group behaviors that have not yet been demonstrated in self-organized robots. For instance, it is not yet
 291 understood how to have robots autonomously identify when they should make a collective decision (Khaluf
 292 et al., 2019). These fundamentals of group-level autonomy, which social animals such as humans exhibit
 293 effortlessly and consistently, might possibly be based on, or even depend on, such unexpected features
 294 as shared biases. Our perspective is that research that investigates the transfer of such social traits from
 295 humans to robots can help us to identify and understand the basic elements needed to build artificial social
 296 cognition.

297 Artificial restrictions in embodied experiments are unlikely to reveal how humans would behave in natural
 298 conditions, but there is existing evidence that such restrictions indeed have the potential to reveal aspects
 299 of embodied human social behavior that would be transferable to robots. For example, when realistic
 300 social cues such as gaze and facial expressions are inhibited, humans have been shown to focus on other
 301 communication channels, such as implicit movement-based communication (Roth et al., 2016).

302 If eventually achieved, the creation of social cognition among robots would open many further research
 303 questions. For instance, there are human collective intentions that go beyond the humans that are

304 immediately present (Tomasello et al., 2005)—if robots have advanced social cognition abilities, how
305 should different social groups of robots interact with each other, whether physically or remotely? As
306 another example, intrinsic motivation or curiosity-driven learning could be investigated to motivate agents
307 to explore the complex internal states that make up another agent, perhaps constituting a rudimentary
308 theory of an artificial mind. Or, perhaps robots could be intrinsically motivated to autonomously develop
309 completely new forms of artificial social cognition that do not resemble those already seen in humans or
310 social animals.

CONFLICT OF INTEREST STATEMENT

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

AUTHOR CONTRIBUTIONS

313 All the authors contributed the ideas and concepts presented in the paper. NC and MKH wrote the first draft
314 of the manuscript. All authors contributed to manuscript revision and read and approved the submitted
315 version.

FUNDING

316 This work was supported by the program of Concerted Research Actions (ARC) of the Université libre de
317 Bruxelles.

ACKNOWLEDGMENTS

318 M. K. Heinrich, A. Cleeremans and M. Dorigo acknowledge support from the F.R.S.-FNRS, of which they
319 are, respectively, postdoctoral researcher and research directors.

REFERENCES

- 320 Albrecht, S. V. and Stone, P. (2018). Autonomous agents modelling other agents: A comprehensive survey
321 and open problems. *Artificial Intelligence* 258, 66–95. doi:10.1016/j.artint.2018.01.002
- 322 Allwright, M., Zhu, W., and Dorigo, M. (2019). An open-source multi-robot construction system.
323 *HardwareX* 5, e00050. doi:10.1016/j.ohx.2018.e00050
- 324 Almaatouq, A., Noriega-Campero, A., Alotaibi, A., Krafft, P. M., Moussaid, M., and Pentland, A. (2020).
325 Adaptive social networks promote the wisdom of crowds. *Proceedings of the National Academy of*
326 *Sciences* 117, 11379–11386. doi:10.1073/pnas.1917687117
- 327 Baker, C. L., Saxe, R., and Tenenbaum, J. B. (2009). Action understanding as inverse planning. *Cognition*
328 113, 329–349. doi:10.1016/j.cognition.2009.07.005
- 329 Bard, N., Foerster, J. N., Chandar, S., Burch, N., Lanctot, M., Song, H. F., et al. (2020). The hanabi
330 challenge: A new frontier for ai research. *Artificial Intelligence* 280, 103216
- 331 Bolotta, S. and Dumas, G. (2022). Social neuro AI: Social interaction as the “dark matter” of AI. *Frontiers*
332 *in Computer Science* 4. doi:10.3389/fcomp.2022.846440
- 333 Botvinick, M. M., Braver, T. S., Barch, D. M., Carter, C. S., and Cohen, J. D. (2001). Conflict monitoring
334 and cognitive control. *Psychological Review* 108, 624–652. doi:10.1037/0033-295x.108.3.624

- 335 Buzsáki, G. (2019). *The Brain from Inside Out* (Oxford University Press). doi:10.1093/oso/9780190905385.
336 001.0001
- 337 Cangelosi, A. and Asada, M. (2022). *Cognitive robotics* (MIT Press)
- 338 Coucke, N., Heinrich, M. K., Cleeremans, A., and Dorigo, M. (2020). Hugos: A multi-user virtual
339 environment for studying human–human swarm intelligence. In *International conference on swarm*
340 *intelligence* (Springer), 161–175
- 341 Coucke, N., Heinrich, M. K., Cleeremans, A., and Dorigo, M. (2021). HuGoS: a virtual environment
342 for studying collective human behavior from a swarm intelligence perspective. *Swarm Intelligence*
343 doi:10.1007/s11721-021-00199-1
- 344 Dale, R., Fusaroli, R., Duran, N. D., and Richardson, D. C. (2013). The self organization of human
345 interaction. *Psychology of Learning and Motivation* 59, 43–95
- 346 Daniel, C., Kroemer, O., Viering, M., Metz, J., and Peters, J. (2015). Active reward learning with a novel
347 acquisition function. *Autonomous Robots* 39, 389–405
- 348 Dautenhahn, K. (2007). Socially intelligent robots: dimensions of human–robot interaction. *Philosophical*
349 *Transactions of the Royal Society B: Biological Sciences* 362, 679–704. doi:10.1098/rstb.2006.2004
- 350 De Jaegher, H., Di Paolo, E., and Gallagher, S. (2010). Can social interaction constitute social cognition?
351 *Trends in cognitive sciences* 14, 441–447
- 352 Dorigo, M., Theraulaz, G., and Trianni, V. (2020). Reflections on the future of swarm robotics. *Science*
353 *Robotics* 5, eabe4385
- 354 Dorigo, M., Theraulaz, G., and Trianni, V. (2021). Swarm robotics: past, present, and future [point of
355 view]. *Proceedings of the IEEE* 109, 1152–1165
- 356 Dumas, G. and Fairhurst, M. T. (2021). Reciprocity and alignment: quantifying coupling in dynamic
357 interactions. *Royal Society Open Science* 8. doi:10.1098/rsos.210138
- 358 Friston, K. and Frith, C. (2015). A duet for one. *Consciousness and Cognition* 36, 390–405. doi:10.1016/j.
359 concog.2014.12.003
- 360 Frith, C. D. (2008). Social cognition. *Philosophical Transactions of the Royal Society B: Biological*
361 *Sciences* 363, 2033–2039
- 362 Frith, C. D. and Frith, U. (2012). Mechanisms of social cognition. *Annual Review of Psychology* 63,
363 287–313. doi:10.1146/annurev-psych-120710-100449
- 364 Frohnwieser, A., Murray, J. C., Pike, T. W., and Wilkinson, A. (2016). Using robots to understand animal
365 cognition. *Journal of the experimental analysis of behavior* 105, 14–22
- 366 Gordon, J., Maselli, A., Lancia, G. L., Thiery, T., Cisek, P., and Pezzulo, G. (2021). The road towards
367 understanding embodied decisions. *Neuroscience & Biobehavioral Reviews* 131, 722–736. doi:10.1016/
368 j.neubiorev.2021.09.034
- 369 Haggard, P. and Chambon, V. (2012). Sense of agency. *Current Biology* 22, R390–R392. doi:10.1016/j.
370 cub.2012.02.040
- 371 He, H., Boyd-Graber, J., Kwok, K., and Daumé III, H. (2016). Opponent modeling in deep reinforcement
372 learning. In *International conference on machine learning* (PMLR), 1804–1813
- 373 Heinrich, M. K., Wahby, M., Dorigo, M., and Hamann, H. (2022). Swarm robotics. In *Cognitive robotics*,
374 eds. A. Cangelosi and M. Asada (MIT Press). 77–98
- 375 Henschel, A., Hortensius, R., and Cross, E. S. (2020). Social cognition in the age of human–robot
376 interaction. *Trends in Neurosciences* 43, 373–384
- 377 Hertwig, R. and Herzog, S. M. (2009). Fast and frugal heuristics: Tools of social rationality. *Social*
378 *Cognition* 27, 661–698. doi:10.1521/soco.2009.27.5.661

- 379 Kameda, T., Toyokawa, W., and Tindale, R. S. (2022). Information aggregation and collective
380 intelligence beyond the wisdom of crowds. *Nature Reviews Psychology* 1, 345–357. doi:10.1038/
381 s44159-022-00054-y
- 382 Kengyel, D., Hamann, H., Zahadat, P., Radspieler, G., Wotawa, F., and Schmickl, T. (2015). Potential of
383 heterogeneity in collective behaviors: A case study on heterogeneous swarms. In *PRIMA 2015: Principles
384 and Practice of Multi-Agent Systems*, eds. Q. Chen, P. Torrioni, S. Villata, J. Hsu, and A. Omicini (Cham:
385 Springer International Publishing), 201–217. doi:doi.org/10.1007/978-3-319-25524-8_13
- 386 Khaluf, Y., Simoens, P., and Hamann, H. (2019). The neglected pieces of designing collective decision-
387 making processes. *Frontiers in Robotics and AI* 6, 16
- 388 Lokesh, R., Sullivan, S., Calalo, J. A., Roth, A., Swanik, B., Carter, M. J., et al. (2022). Humans
389 utilize sensory evidence of others' intended action to make online decisions. *Scientific Reports* 12.
390 doi:10.1038/s41598-022-12662-y
- 391 Mathews, N., Christensen, A. L., O'Grady, R., Mondada, F., and Dorigo, M. (2017). Mergeable nervous
392 systems for robots. *Nature communications* 8, 1–7
- 393 Moussaid, M., Helbing, D., and Theraulaz, G. (2011). How simple rules determine pedestrian behavior
394 and crowd disasters. *Proceedings of the National Academy of Sciences* 108, 6884–6888. doi:10.1073/
395 pnas.1016507108
- 396 Murata, S., Yamashita, Y., Arie, H., Ogata, T., Sugano, S., and Tani, J. (2015). Learning to perceive the
397 world as probabilistic or deterministic via interaction with others: A neuro-robotics experiment. *IEEE
398 transactions on neural networks and learning systems* 28, 830–848
- 399 Newman-Norlund, R. D., van Schie, H. T., van Zuijlen, A. M. J., and Bekkering, H. (2007). The
400 mirror neuron system is more active during complementary compared with imitative action. *Nature
401 Neuroscience* 10, 817–818. doi:10.1038/nn1911
- 402 Nouyan, S., Groß, R., Bonani, M., Mondada, F., and Dorigo, M. (2009). Teamwork in self-organized robot
403 colonies. *IEEE Transactions on Evolutionary Computation* 13, 695–711. doi:dx.doi.org/10.1109/TEVC.
404 2008.2011746
- 405 Olsson, A., Knapska, E., and Lindström, B. (2020). The neural and computational systems of social
406 learning. *Nature Reviews Neuroscience* 21, 197–212. doi:10.1038/s41583-020-0276-4
- 407 Patel, D., Fleming, S. M., and Kilner, J. M. (2012). Inferring subjective states through the observation of
408 actions. *Proc. R. Soc. B* 279, 4853–4860. doi:10.1098/rspb.2012.1847
- 409 Petersen, K. H., Napp, N., Stuart-Smith, R., Rus, D., and Kovac, M. (2019). A review of collective robotic
410 construction. *Science Robotics* 4. doi:10.1126/scirobotics.aau8479
- 411 Pezzulo, G. and Dindo, H. (2011). What should i do next? using shared representations to solve interaction
412 problems. *Experimental Brain Research* 211, 613–630. doi:10.1007/s00221-011-2712-1
- 413 Rabinowitz, N., Perbet, F., Song, F., Zhang, C., Eslami, S. M. A., and Botvinick, M. (2018). Machine
414 theory of mind. In *Proceedings of the 35th International Conference on Machine Learning*, eds. J. Dy
415 and A. Krause (PMLR), vol. 80 of *Proceedings of Machine Learning Research*, 4218–4227
- 416 Raileanu, R., Denton, E., Szlam, A., and Fergus, R. (2018). Modeling others using oneself in multi-agent
417 reinforcement learning. In *International conference on machine learning* (PMLR), 4257–4266
- 418 Ratcliff, R. and McKoon, G. (2008). The diffusion decision model: Theory and data for two-choice
419 decision tasks. *Neural Computation* 20, 873–922. doi:10.1162/neco.2008.12-06-420
- 420 Rizzolatti, G. and Craighero, L. (2004). The mirror-neuron system. *Annual Review of Neuroscience* 27,
421 169–192. doi:10.1146/annurev.neuro.27.070203.144230

- 422 Roth, D., Lugrin, J.-L., Galakhov, D., Hofmann, A., Bente, G., Latoschik, M. E., et al. (2016). Avatar
423 realism and social interaction quality in virtual reality. In *2016 IEEE Virtual Reality (VR)* (IEEE),
424 277–278. doi:10.1109/vr.2016.7504761
- 425 Rubenstein, M., Cornejo, A., and Nagpal, R. (2014). Programmable self-assembly in a thousand-robot
426 swarm. *Science* 345, 795–799
- 427 Saxe, R. (2005). Against simulation: the argument from error. *Trends in Cognitive Sciences* 9, 174–179.
428 doi:10.1016/j.tics.2005.01.012
- 429 Schilbach, L., Timmermans, B., Reddy, V., Costall, A., Bente, G., Schlicht, T., et al. (2013).
430 Toward a second-person neuroscience. *Behavioral and Brain Sciences* 36, 393–414. doi:10.1017/
431 s0140525x12000660
- 432 Sebanz, N., Bekkering, H., and Knoblich, G. (2006). Joint action: bodies and minds moving together.
433 *Trends in Cognitive Sciences* 10, 70–76. doi:10.1016/j.tics.2005.12.009
- 434 Seymour, B. and Dolan, R. (2008). Emotion, decision making, and the amygdala. *Neuron* 58, 662–671.
435 doi:10.1016/j.neuron.2008.05.020
- 436 Silver, C. A., Tatler, B. W., Chakravarthi, R., and Timmermans, B. (2020). Social agency as a continuum.
437 *Psychonomic Bulletin & Review* 28, 434–453. doi:10.3758/s13423-020-01845-1
- 438 Tomasello, M., Carpenter, M., Call, J., Behne, T., and Moll, H. (2005). Understanding and sharing
439 intentions: The origins of cultural cognition. *Behavioral and Brain Sciences* 28, 675–691. doi:10.1017/
440 s0140525x05000129
- 441 Tump, A. N., Pleskac, T. J., and Kurvers, R. H. J. M. (2020). Wise or mad crowds? the cognitive
442 mechanisms underlying information cascades. *Science Advances* 6. doi:10.1126/sciadv.abb0266
- 443 Ullman, T., Baker, C., Macindoe, O., Evans, O., Goodman, N., and Tenenbaum, J. (2009). Help or hinder:
444 Bayesian models of social goal inference. In *Advances in Neural Information Processing Systems*, eds.
445 Y. Bengio, D. Schuurmans, J. Lafferty, C. Williams, and A. Culotta (Curran Associates, Inc.), vol. 22,
446 1874–1882
- 447 Valentini, G., Ferrante, E., and Dorigo, M. (2017). The best-of-n problem in robot swarms: Formalization,
448 state of the art, and novel perspectives. *Frontiers in Robotics and AI* 4. doi:10.3389/frobt.2017.00009
- 449 Valentini, G., Ferrante, E., Hamann, H., and Dorigo, M. (2016). Collective decision with 100 kilobots:
450 Speed versus accuracy in binary discrimination problems. *Autonomous agents and multi-agent systems*
451 30, 553–580
- 452 Van Vugt, M. (2006). Evolutionary origins of leadership and followership. *Personality and Social*
453 *Psychology Review* 10, 354–371. doi:10.1207/s15327957pspr1004_5
- 454 Vernon, D. (2014). *Artificial cognitive systems: A primer* (MIT Press)
- 455 Wilson, M. (2002). Six views of embodied cognition. *Psychonomic Bulletin & Review* 9, 625–636.
456 doi:10.3758/bf03196322
- 457 Wykowska, A., Chaminade, T., and Cheng, G. (2016). Embodied artificial agents for understanding human
458 social cognition. *Philosophical Transactions of the Royal Society B: Biological Sciences* 371, 20150375
- 459 Zhu, W., Allwright, M., Heinrich, M. K., Oğuz, S., Christensen, A. L., and Dorigo, M. (2020). Formation
460 control of uavs and mobile robots using self-organized communication topologies. In *International*
461 *conference on swarm intelligence* (Springer), 306–314