
SOCIAL NEURO AI: SOCIAL INTERACTION AS THE "DARK MATTER" OF AI

A PREPRINT

✉ **Samuele Bolotta**

Department of Computer Science and Operations Research
Université de Montréal
Montreal, QC, H2S 3H1
samuele.bolotta@ppsp.team

✉ **Guillaume Dumas**

Mila - Quebec Artificial Intelligence Institute
CHU Sainte-Justine Research Center, Department of Psychiatry
Université de Montréal
Montreal, QC, H2S 3H1
guillaume.dumas@ppsp.team

January 4, 2022

ABSTRACT

We are making the case that empirical results from social psychology and social neuroscience along with the framework of dynamics can be of inspiration to the development of more intelligent artificial agents. We specifically argue that the complex human cognitive architecture owes a large portion of its expressive power to its ability to engage in social and cultural learning. In the first section, we aim at demonstrating that social learning plays a key role in the development of intelligence. We do so by discussing social and cultural learning theories and investigating the abilities that various animals have at learning from others; we also explore findings from social neuroscience that examine human brains during social interaction and learning. Then, we discuss three proposed lines of research that fall under the umbrella of Social NeuroAI and can contribute to developing socially intelligent embodied agents in complex environments. First, neuroscientific theories of cognitive architecture, such as the global workspace theory and the attention schema theory, can enhance biological plausibility and help us understand how we could bridge individual and social theories of intelligence. Second, intelligence occurs in time as opposed to over time, and this is naturally incorporated by the powerful framework offered by dynamics. Third, social embodiment has been demonstrated to provide social interactions between virtual agents and humans with a more sophisticated array of communicative signals. To conclude, we provide a new perspective on the field of multiagent robot systems, exploring how it can advance by following the aforementioned three axes.

1 The importance of social learning

Social learning categories Human cognitive functions such as theory of mind and explicit metacognition are not genetically programmed, but rather constructed as “cognitive gadgets” during development through social interaction Heyes [2018]. Since their birth, social animals use their conspecifics as vehicles for gathering information that can potentially help them respond efficiently to challenges in the environment, avoiding harm and maximizing rewards Kendal et al. [2018]. Learning adaptive information from others results in better regulation of task performance, especially by gaining fitness benefits and in avoiding some of the costs associated with asocial, trial-and-error learning, such as time loss, energy loss, and exposure to predation Clark and Dumas [2016]. Importantly, cultural inheritance permeates a broad array of behavioural domains, including migratory pathways, foraging techniques, nesting sites and mates Whiten [2021]. The spread of such information across generations gives social learning a unique role in the evolution of culture and therefore makes it a crucial candidate to investigate the biological bases of human cognition. A large body of evidence suggests that culture is not specific to humans, with surprising discoveries involving insects. Recent reviews have identified four main categories of social learning that differ in what is socially learnt and in the cognitive skills that are required Hopppitt and Laland [2008], Whiten [2021] (Figure 1). These categories have been developed through the approach of behaviourism. While we acknowledge that there is more to social learning than mere behaviour (the affective and emotional dimensions are equally crucial Gruber et al. [2021]), we keep it as the focus of this short article because it is an empirically solid starting point with clarified mecha-

nisms. At the most elementary level, enhancement consists of an agent observing a model that focuses on particular objects or locations and consequently adopting the same focus Heyes [1994], Thorpe [1963]. For example, it was demonstrated that bees outside the nest land more often on flowers that they had seen preferred by the model Worden and Papaj [2005]. This skill requires social agents to perform basic associative learning in relation to other agents’ observed actions; it is likely to be the most widespread form of social learning across the animal kingdom. A more complex form of social learning consists of observational conditioning, which exposes a social agent to a relationship between stimuli Heyes [1994]; this exposure causes a change in the agent. For example, the observation of experienced demonstrators facilitated the opening of hickory nuts by red squirrels, relative to trial-and-error learning Weigl and Hanson [1980]. This is therefore a mechanism through which agents learn the value of a stimulus from the interaction with other agents. Yet a more complex form of social learning consists of affordance learning, which allows a social agent to learn the operating characteristics of objects or environments by observing the behaviour of other agents Whiten [2021]. For example, pigeons that saw a demonstrator push a sliding screen for food made a higher proportion of pushes than observers in control conditions, thus exhibiting affordance learning Klein and Zentall [2003]. In other words, the animals learn the effects that a certain action has on the environment. Finally, at the most complex level, copying another individual can take the shape of pure imitation, where every detail is copied, or emulation, where only a few elements are copied Byrne [2002]. For example, most chimpanzees mastered a new technique for obtaining food when they were under the influence of a trained expert, whereas none did so in a population lacking an expert Whiten [2005].

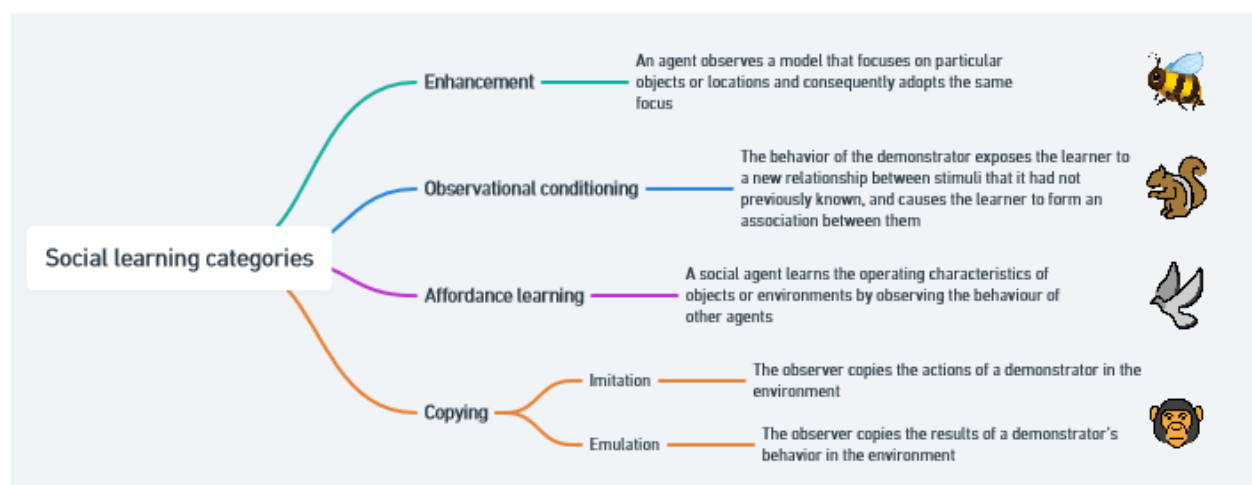


Figure 1: Social learning categories. Figure inspired by Whiten [2005].

Social learning strategies Crucially, while social learning is widespread, using it indiscriminately is rarely beneficial. This suggests that individuals should be selective in what, when, and whom they copy, by following 'social learning strategies' (SLSs; Kendal et al. [2018]). Several SLSs might be used by the same population and even by the same individual. The aforementioned categories of social learning have been shown to be refined by modulating biases that can strengthen their adaptive power Kendal et al. [2018]. For example, an important SLS is copying when asocial learning would be costly; research has shown that, when task difficulty increases, various animals are more likely to use social information. Individuals also prefer using social information when they are uncertain about a task; high-fidelity copying is observed among children who lack relevant personal information Wood et al. [2013]. In general, other state-based SLSs can affect the decision to use social information, such as age, social rank, and reproductive state of the learner; for example, low- and mid-ranking chimpanzees are more likely to use social information than high-ranking individuals Kendal et al. [2015]. Model-based biases are another crucial category; for example, children prefer to copy prestigious individuals, where status is evidenced by their older age, popularity and social dominance Flynn and Whiten [2012]. Multiple evidence also suggests that a conformist transmission bias exists, whereby the behaviour of the majority of individuals is more likely to be adopted by others Kendal et al. [2018].

Social learning in neuroscience Despite the progress made in social neuroscience and in developmental psychology, only in the last decade, serious efforts have started focusing on the neural mechanisms of social interaction, which were seen as the "dark matter" of social neuroscience Schilbach et al. [2013]; recently, a framework for computational social neuroscience has been proposed, in an attempt to naturalize social interaction Tognoli et al. [2018]. At the intra-brain level, it was demonstrated that social interaction is categorically different from social perception and that the brain exhibits different activity patterns depending on the role of the subject and on the context in which the interaction is unfolding Dumas et al. [2012]. At the inter-brain level, functional Magnetic Resonance Imaging (fMRI) or Electroencephalography (EEG) recordings of multiple brains (i.e. hyperscanning) have allowed to demonstrate inter-brain synchronization during social interaction - specifically, while subjects were engaged in spontaneous imitation of hand movements Dumas et al. [2010]. Interestingly, the increase in coupling strength between brain signals was also shown to be present during a two-person turn-taking verbal exchange with no visual contact, in both a native or a foreign language context Pérez et al. [2019]. Inter-brain synchronization is also modulated by the type of task and by the familiarity between subjects Djalovski et al. [2021]. Overall, this shows that, beyond their individual cognition, humans are also coupled in the social dimension. Interestingly, the field of computational social neuroscience has also focused on explaining the functional meaning of such correlations between inter-

brain synchronization and behavioural synchrony. A biophysical model showed that the similarity of endogenous dynamics and the similarity of anatomical structure might facilitate interindividual synchronization and explain our propensity to generate interindividual coupling via perception and actions Dumas et al. [2012]. Moreover, the topology of the anatomy of the individual brain seems to enhance the resonance between brains; this suggests that beyond facilitating integration of information within the brain, the human connectome tends to facilitate sensorimotor resonance between individuals Dumas et al. [2012]. This evidence clearly points to the crucial importance that social connections have in humans.

Social learning and language development Regarding language development in humans, cognitive and structural accounts of language development have often conceptualized linguistic abilities as static and formal sets of knowledge structures, ignoring the contextual nature of language. However, good communication must be tailored to the characteristics of the listener and of the context - language can also be explained as a social construct Whitehurst [1978]. For example, evidence shows that the language outcome of children with cochlear implants is heavily influenced by parental linguistic input during the first years after cochlear implantation Holzinger et al. [2020]. In terms of specific social learning variables, imitation has also been shown to play a major role in boosting language development, usually in the form of selective imitation Whitehurst et al. [1974], Whitehurst and Vasta [1975]. Moreover, in autistic kids, social learning variables such as joint attention, imitative imitation, and deferred imitation have been shown to be the best predictors of language ability and rate of communication development Toth et al. [2006]. These results clearly suggest that social learning skills have an influence on language acquisition in humans.

2 Steps towards Social Neuro AI

How could social learning be useful for AI? As AI will play a crucial role in our daily lives, one of the main challenges is building autonomous agents capable of participating in cooperative social interactions with humans. In the previous sections, we have provided convincing evidence that interpersonal intelligence enhances intrapersonal intelligence through the mechanisms and biases of social learning. It is a crucial hallmark of many species and it manifests itself across different behavioural domains; without it, animals would lose the possibility to quickly acquire valuable information from their conspecifics and therefore lose fitness benefits. At the same time, social learning cannot be used indiscriminately and a broad array of modulating biases exists to strengthen its adaptive power. Recent efforts in computational social neuroscience have paved the way for a naturalization of social interactions, showing that the connectome seems to facilitate resonance between brains; beyond their individual cognition, humans are also dynamically coupled in the social realm. Regarding artificial intelligence, multi-agent reinforcement

learning (MARL) is the best subfield to investigate the interactions between multiple agents. Such interactions can be of three types: cooperative games (all agents working for the same goal), competitive games (all agents competing against each other), and mixed motive games (a mix of cooperative and competitive interactions). At each timestep t , each agent is attempting to maximize its own reward by learning a policy that optimizes the total expected discounted future reward. We refer the reader to high-quality reviews that have been written on MARL Wong et al. [2021], Nguyen et al. [2020], Hernandez-Leal et al. [2019]. Here, we highlight that, among others, low sample efficiency is one of the greatest challenges for MARL, as millions of interactions with the environment are usually needed for agents to learn. Moreover, multi-agent joint action space increases exponentially with the number of agents, leading to problems that are often intractable. In the last few years, part of the AI community has already started demonstrating that these problems can be alleviated by mechanisms that allow for social learning Jaques [2019], Ndousse [2021], Lee et al. [2021]. More in general, concepts from complex systems such as self-organization, emergent behavior, swarm optimization and cellular systems suggest that collective intelligence could produce more robust and flexible solutions in AI, with higher sample efficiency and higher generalization Ha and Tang [2021]. In the following sections, we argue that to exploit all benefits that social learning can offer, more focus on biological plausibility, social embodiment and temporal dynamics is needed. These three approaches are all necessary, but not sufficient by themselves, to generate social learning.

2.1 Biological plausibility.

The social learning skills and biases that we have shown so far are boosted in humans by their advanced cognitive architecture Whiten [2021]. Equipping artificial agents with complex social learning abilities will therefore require more complex architectures that can handle a great variety of information efficiently. While the human unconscious brain aligns well with the current successful applications of deep learning, the conscious brain involves higher-order cognitive abilities that perform much more complex computations than what deep learning can currently do Bengio [2019]. More specifically, "unconsciousness" is where most of our intelligence lies and involves unconscious abilities related to view-invariance, meaning extraction, control, decision-making and learning; "i-consciousness" is the part of human consciousness that is focused on integrating all available evidence to converge toward a single decision; "m-consciousness" is the part of human consciousness that is focused on reflexively representing oneself, utilizing error detection, meta-memory and reality monitoring Graziano [2017]. Notably, recent efforts in the deep learning community have indeed focused on building advanced cognitive architectures that are inspired from neuroscience. In particular, the global workspace theory (GWT) is the most widely accepted theory of consciousness, and it postulates that when a piece of information

is selected by attention, it may non-linearly achieve "ignition", enter the global workspace (GLW) and be shared across specialized cortical modules, therefore becoming conscious Baars [1993]; Dehaene et al. [1998]. The use of such a communication channel in the context of deep learning was explored for modelling the structure of complex environments. This architecture was demonstrated to encourage specialization and compositionality and to facilitate the synchronization of otherwise independent specialists Goyal et al. [2021]. Moreover, inductive biases inspired by higher-order cognitive functions in humans have been shown to improve out-of-distribution generalization. However, the GWT is correct, but incomplete Zhao et al. [2021]. One important principle in control engineering is that a good controller contains a model of the item being controlled Conant and Ross Ashby [1970]. The attention schema theory (AST) tackles this problem by adding an attention schema to the GLW. The performance of an artificial agent in solving a simple sensorimotor task is greatly enhanced by an attention schema, but its performance is greatly reduced when the schema is not available Wilterson and Graziano [2021]. Specifically, the proposal is that the brain constructs not only a model of the physical body but also a coherent, rich, and descriptive model of attention. The body schema contains layers of valuable information that help control and predict stable and dynamic properties of the body; in a similar fashion, the attention schema helps control and predict attention. One cannot understand how the brain controls the body without understanding the body schema, and in a similar way one cannot understand how the brain controls its limited resources without understanding the attention schema Graziano [2017]. Therefore, the study of consciousness in artificial intelligence is not a mere pursuit of metaphysical mystery; from an engineering perspective, without understanding subjective awareness, it might not be possible to build artificial agents that intelligently control and deploy their limited processing resources. It has also been argued that, without an attention schema, it might be impossible to build artificial agents that are socially intelligent. This idea stems from the evidence that points at an overlap of social cognition functions with awareness and attention functions in the right temporo-parietal junction of the human brain Mitchell [2008]. It was then proposed that an attention schema might also be used for social cognition, giving rise to an overlap between modelling one's own attention and modelling others' attention. In other words, when we attribute to other people an awareness of their surroundings, we are constructing a simplified model of their attention - a schema of others' attention Graziano and Kastner [2011]. Overall, this section proposes that we make a transition from Good Old-Fashioned Artificial Intelligence to Neuro-AI, so as to draw inspiration from one structure we know is capable of complex intelligence: the human brain (Figure 2).

2.2 Temporal dynamics.

In nature, complex systems are composed of simple components that self-organize in time, producing ultimately

emergent behaviours that depend on the dynamical interactions between the components. Such focus on coordination dynamics can help in shifting the perspective from representation-centered to self-organizing agents Brooks [1991]. The former view has been one predominant way of thinking about autonomous systems that exhibit intelligent behaviour: such autonomous agents use their sensors to extract information about the world they operate in and use it to construct an internal model of the world and therefore rationally perform optimal decision making in pursuit of some goal. In other words, autonomous agents are information processing systems and their environment can be abstracted away as the source of answers to questions raised by the ongoing agents' needs. Importantly, according to this view, sensorimotor connections of the agents to the environment are still relevant to understand their behaviour, but there is no focus on what such connections involve and how they take place Newell and Simon [1976]. The latter view, in line with Brooks' ideas, shows how the representational approach ignores the nonlinear dynamical aspect of intelligence, that is, the temporal constraints that characterize the interactions between agent and environment. Instead, dynamics is a powerful framework that has been used to describe multiple natural phenomena as an interdependent set of coevolving quantitative variables van Gelder [1998] and a crucial aspect of intelligence is that it occurs in time and not over time. If we abstract away the richness of real time, then we also change the behaviour of the agents Smithers [2018]. In other words, one should indeed focus on the structural complexity and on the algorithmic computation the agents need to carry out, but without abstracting away the dynamical aspects of the agent-environment interactions: such dynamical aspects are pervasive and, therefore, necessary to explain the behaviour of the system van Gelder [1998], Smithers [2018], Barandiaran [2017]. Connected to these ideas, a variational principle of least free energy Friston et al. [2006], also known as active inference, sees dynamical systems as cognitive systems that instantiate in a generative model one main goal: minimizing on average their surprise, through perception and action Ramstead et al. [2020]. According to this view, generative models do not encode exploitable and symbolic structural information about the world; they are control systems that are expressed in embodied activity and utilize information encoded in the approximate posterior belief Ramstead et al. [2020]. Active inference models are still very discrete in their architectures, especially regarding high-level behavioral aspects, but they may be a good class of models to raise the tension between computation and implementation (Figure 2).

2.3 Social embodiment.

There has been a resurgence of enactivism in cognitive neuroscience over the past decade, emphasizing the circular causality induced by the notion that the environment is acting upon the individual and the individual is acting upon the environment. To understand how the brain works, then one has to acknowledge that it is embodied Clark [2013], Hohwy [2013]. Evidence for this shows that embodied in-

telligence in human children arises from the interaction of the child with the environment through a sensory body that is capable of recognizing the statistical properties of such interaction Smith and Gasser [2005]. Moreover, higher primates interpret each other as psychological subjects based on their bodily presence; social embodiment is the idea that the embodiment of a socially interactive agent plays a significant role in social interactions. It refers to "states of the body, such as postures, arm movements, and facial expressions, that arise during social interaction and play central roles in social information processing." Thompson and Varela [2001], Barsalou et al. [2003]. This includes internal and external structures, sensors, and motors that allow them to interact actively with the world. At a high level, sensorimotor capabilities in the avatar and robots are meant to model their role in biological beings: the agent now has limitations in the ways they can sense, manipulate, and navigate its environments. Importantly, these limitations are closely tied to the agent's function Deng et al. [2019]. The idea of social embodiment in artificial agents is supported by evidence of improvements in the interactions between embodied agents and humans Zhang et al. [2016]. Studies have shown positive effects of physical embodiment on the feeling of an agent's social presence, the evaluation of the agent, the assessment of public evaluation of the agent, and the evaluation of the interaction with the agent Kose-Bagci et al. [2009], Gupta et al. [2021]. In robots, social presence is a key component in the success of social interactions and it can be defined as the combination of seven abilities that enhance a robot's social skills: 1. Express emotion 2. Communicate with high-level dialogue 3. Learn/recognize models of other agents 4. Establish/maintain social relationships 5. Use natural cues 6. Exhibit distinctive personality and character 7. Learn/develop social competencies Lee [2006]. Social embodiment thus equips artificial agents with a more articulated and richer repertoire of expressions, ameliorating the interactions with it Jaques [2019]. For instance, in human-robot interaction, a gripper is not limited to its role in the manipulation of objects. Rather, it opens a broad array of movements that can enhance the communicative skills of the robot and, consequently, the quality of its possible interactions Deng et al. [2019]. The embodied agent is therefore the best model of the aspects of the world relevant to its surviving and thriving, through performing situationally appropriate actions Ramstead et al. [2020] (Figure 2).

3 Conclusion

At the crossroads of robotics, computer science, and psychology, one of the main challenges for humans is to build autonomous agents capable of participating in cooperative social interactions. This is important not only because AI will play a crucial role in our daily life, but also because, as demonstrated by results in social neuroscience and evolutionary psychology, intrapersonal intelligence is tightly connected with interpersonal intelligence, especially in humans Dumas et al. [2014a]. In this opinion article, we

have attempted to unify the lines of research that, at the moment, are separated from each other; in particular, we have proposed three research directions that are expected to enhance efficient exchange of information between agents and, as a consequence, individual intelligence (especially in out-of-distribution generalization: OOD). This would contribute to creating agents that not only do have human-like OOD skills, but are also able to exhibit such skills in extremely complex and realistic environments Dennis et al. [2021], while interacting with other embodied agents and with humans. In parallel, then, it will be crucial to scale up the realism of what the agents perceive in their social context, going from simple environments like GridWorld to more complex ones powered by video-game engines and, finally, to extremely realistic environments, like the one offered by the MetaHuman Creator of Unreal Engine. These

advancements will hopefully result in more socially intelligent agents and therefore in more fruitful interactions between humans and virtual agents. It is not by coincidence that Alan Turing chose social interaction as the ultimate test for machine intelligence Turing [1950]. We previously argued that the genuineness of this subjective judgment of humanness can be strengthened by adding a criterion about the implementation and mechanisms underlying the social behavior of the machine Dumas et al. [2014b]. Here, we have argued that passing this updated Turing test may require the synergy of the two types of computation introduced by Alan Turing: the discrete symbolic computation at the very basis of contemporary digital technologies Turing et al. [1936] and the dynamical embodied computation at play in morphogenesis Turing [1952].

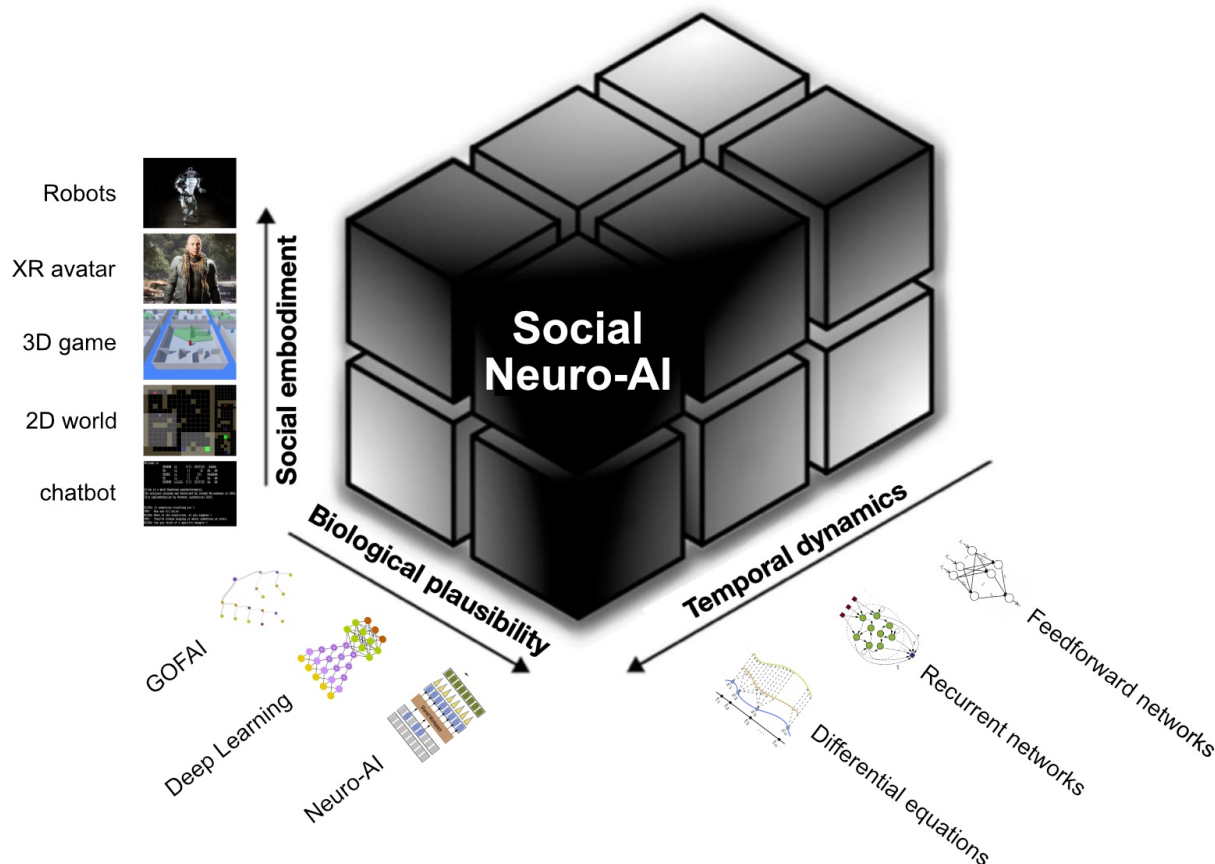


Figure 2: Billions of humans interact daily with algorithms — yet AI is far from human social cognition. We argue that creating such socially aware agents may require "Social Neuro-AI" — a program developing 3 research axes: 1. Biological plausibility 2. Temporal dynamics 3. Social embodiment. Overall, those steps towards socially aware agents will ultimately help in aligned interactions between natural and artificial intelligence. Figure inspired by Schilbach et al. [2013].

References

- Cecilia Heyes. *Cognitive Gadgets*. Harvard University Press, April 2018. ISBN 9780674985155. URL <https://www.degruyter.com/document/doi/10.4159/9780674985155/html>.
- Rachel L. Kendal, Neeltje J. Boogert, Luke Rendell, Kevin N. Laland, Mike Webster, and Patricia L. Jones. Social learning strategies: Bridge-building between fields. *Trends in Cognitive Sciences*, 22(7): 651–665, 2018. doi:10.1016/j.tics.2018.04.003. URL <https://linkinghub.elsevier.com/retrieve/pii/S1364661318300949>.
- Ian Clark and Guillaume Dumas. The regulation of task performance: a trans-disciplinary review. *Frontiers in psychology*, 6:1862, 2016.
- Andrew Whiten. The burgeoning reach of animal culture. *Science*, 372(6537):eabe6514, April 2021. doi:10.1126/science.abe6514. URL <https://www.sciencemag.org/lookup/doi/10.1126/science.abe6514>.
- Will Hoppitt and Kevin N. Laland. *Chapter 3 Social Processes Influencing Learning in Animals: A Review of the Evidence*, volume 38, page 105–165. Elsevier, 2008. ISBN 9780120045389. doi:10.1016/S0065-3454(08)00003-X. URL <https://linkinghub.elsevier.com/retrieve/pii/S006534540800003X>.
- Thibaud Gruber, Marina Bazhydai, Christine Sievers, Fabrice Clément, and Daniel Dukes. The abc of social learning: Affect, behavior, and cognition. *Psychological Review*, Jul 2021. doi:10.1037/rev0000311.
- C. M. Heyes. Social learning in animals: categories and mechanisms. *Biological Reviews of the Cambridge Philosophical Society*, 69(2):207–231, May 1994. doi:10.1111/j.1469-185x.1994.tb01506.x.
- W. H Thorpe. *Learning and instinct in animals*. Methuen, 1963.
- Bradley D Worden and Daniel R Papaj. Flower choice copying in bumblebees. *Biology Letters*, 1(4):504–507, December 2005. doi:10.1098/rsbl.2005.0368. URL <https://royalsocietypublishing.org/doi/10.1098/rsbl.2005.0368>.
- Peter D. Weigl and Elinor V. Hanson. Observational learning and the feeding behavior of the red squirrel *tamiasciurus hudsonicus*: The ontogeny of optimization. *Ecology*, 61(2):213–218, 1980. doi:https://doi.org/10.2307/1935176. URL <https://esajournals.onlinelibrary.wiley.com/doi/abs/10.2307/1935176>.
- Emily D. Klein and Thomas R. Zentall. Imitation and affordance learning by pigeons (*columba livia*). *Journal of Comparative Psychology*, 117(4): 414–419, 2003. doi:10.1037/0735-7036.117.4.414. URL <http://doi.apa.org/getdoi.cfm?doi=10.1037/0735-7036.117.4.414>.
- Richard W. Byrne. *Imitation of novel complex actions: What does the evidence from animals mean?*, volume 31, page 77–105. Academic Press, January 2002. doi:10.1016/S0065-3454(02)80006-7. URL <https://www.sciencedirect.com/science/article/pii/S0065345402800067>.
- Andrew Whiten. The second inheritance system of chimpanzees and humans. *Nature*, 437(7055):52–55, 2005. doi:10.1038/nature04023. URL <http://www.nature.com/articles/nature04023>.
- Lara A. Wood, Rachel L. Kendal, and Emma G. Flynn. Copy me or copy you? the effect of prior experience on social learning. *Cognition*, 127(2):203–213, 2013. doi:10.1016/j.cognition.2013.01.002. URL <https://linkinghub.elsevier.com/retrieve/pii/S0010027713000103>.
- Rachel Kendal, Lydia M. Hopper, Andrew Whiten, Sarah F. Brosnan, Susan P. Lambeth, Steven J. Schapiro, and Will Hoppitt. Chimpanzees copy dominant and knowledgeable individuals: implications for cultural diversity. *Evolution and Human Behavior*, 36(1):65–72, 2015. doi:10.1016/j.evolhumbehav.2014.09.002. URL <https://linkinghub.elsevier.com/retrieve/pii/S109051381400110X>.
- Emma Flynn and Andrew Whiten. Experimental “microcultures” in young children: Identifying biographic, cognitive, and social predictors of information transmission: Identifying predictors of information transmission. *Child Development*, 83(3):911–925, 2012. doi:10.1111/j.1467-8624.2012.01747.x. URL <http://doi.wiley.com/10.1111/j.1467-8624.2012.01747.x>.
- Leonhard Schilbach, Bert Timmermans, Vasudevi Reddy, Alan Costall, Gary Bente, Tobias Schlicht, and Kai Vogeley. Toward a second-person neuroscience. *Behavioral and Brain Sciences*, 36(4):393–414, 2013. doi:10.1017/S0140525X12000660. URL https://www.cambridge.org/core/product/identifier/S0140525X12000660/type/journal_article.
- Emmanuelle Tognoli, Guillaume Dumas, and J. A. Scott Kelso. A roadmap to computational social neuroscience. *Cognitive Neurodynamics*, 12(1):135–140, 2018. doi:10.1007/s11571-017-9462-0. URL <http://link.springer.com/10.1007/s11571-017-9462-0>.
- Guillaume Dumas, Mario Chavez, Jacqueline Nadel, and Jacques Martinerie. Anatomical connectivity influences both intra- and inter-brain synchronizations. *PLoS ONE*, 7(5):e36414, May 2012. doi:10.1371/journal.pone.0036414. URL <https://dx.plos.org/10.1371/journal.pone.0036414>.
- Guillaume Dumas, Jacqueline Nadel, Robert Soussignan, Jacques Martinerie, and Line Garnero. Inter-brain synchronization during social interaction. *PLoS ONE*, 5(8):e12166, August 2010. doi:10.1371/journal.pone.0012166. URL <https://dx.plos.org/10.1371/journal.pone.0012166>.

- Alejandro Pérez, Guillaume Dumas, Melek Karadag, and Jon Andoni Duñabeitia. Differential brain-to-brain entrainment while speaking and listening in native and foreign languages. *Cortex*, 111:303–315, 2019. doi:10.1016/j.cortex.2018.11.026. URL <https://linkinghub.elsevier.com/retrieve/pii/S0010945218304052>.
- Amir Djalovski, Guillaume Dumas, Sivan Kinreich, and Ruth Feldman. Human attachments shape interbrain synchrony toward efficient performance of social goals. *NeuroImage*, 226:117600, 2021. doi:10.1016/j.neuroimage.2020.117600. URL <https://linkinghub.elsevier.com/retrieve/pii/S1053811920310855>.
- Grover J. Whitehurst. The contributions of social learning to language acquisition. *Contemporary Educational Psychology*, 3(1):2–10, Jan 1978. doi:10.1016/0361-476X(78)90002-4. URL <https://www.sciencedirect.com/science/article/pii/0361476X78900024>.
- Daniel Holzinger, Magdalena Dall, Susana Sanduvete-Chaves, David Saldaña, Salvador Chacón-Moscoso, and Johannes Fellingner. The impact of family environment on language development of children with cochlear implants: A systematic review and meta-analysis. *Ear Hearing*, 41(5):1077–1091, 2020. doi:10.1097/AUD.0000000000000852. URL <https://journals.lww.com/10.1097/AUD.0000000000000852>.
- Grover J Whitehurst, Marsha Ironsmith, and Michael Goldfein. Selective imitation of the passive construction through modeling. *Journal of Experimental Child Psychology*, 17(2):288–302, Apr 1974. doi:10.1016/0022-0965(74)90073-3. URL <https://www.sciencedirect.com/science/article/pii/0022096574900733>.
- Grover J. Whitehurst and Ross Vasta. Is language acquired through imitation? *Journal of Psycholinguistic Research*, 4(1):37–59, Jan 1975. doi:10.1007/BF01066989. URL <https://doi.org/10.1007/BF01066989>.
- Karen Toth, Jeffrey Munson, Andrew N. Meltzoff, and Geraldine Dawson. Early predictors of communication development in young children with autism spectrum disorder: Joint attention, imitation, and toy play. *Journal of Autism and Developmental Disorders*, 36(8):993–1005, Nov 2006. doi:10.1007/s10803-006-0137-7. URL <https://doi.org/10.1007/s10803-006-0137-7>.
- Annie Wong, Thomas Bäck, Anna V. Kononova, and Aske Plaat. Multiagent deep reinforcement learning: Challenges and directions towards human-like approaches. *arXiv:2106.15691 [cs]*, Jun 2021. URL <http://arxiv.org/abs/2106.15691>. arXiv: 2106.15691.
- Thanh Thi Nguyen, Ngoc Duy Nguyen, and Saeid Nahavandi. Deep reinforcement learning for multi-agent systems: A review of challenges, solutions and applications. *IEEE Transactions on Cybernetics*, 50(9):3826–3839, 2020. doi:10.1109/TCYB.2020.2977374. URL <http://arxiv.org/abs/1812.11794>. arXiv: 1812.11794.
- Pablo Hernandez-Leal, Bilal Kartal, and Matthew E. Taylor. A survey and critique of multiagent deep reinforcement learning. *Autonomous Agents and Multi-Agent Systems*, 33(6):750–797, 2019. doi:10.1007/s10458-019-09421-1. URL <http://arxiv.org/abs/1810.05587>. arXiv: 1810.05587.
- Natasha Jaques. Social and affective machine learning, 2019. URL <https://www.media.mit.edu/publications/social-and-affective-machine-learning/>.
- Kamal Ndousse. Emergent social learning via multi-agent reinforcement learning, 2021. URL https://scholar.google.com/citations?view_op=view_citation&hl=it&user=8iCb2TwAAAAJ&sortby=pubdate&citation_for_view=8iCb2TwAAAAJ:j3f4tGmQtD8C.
- Dennis Lee, Natasha Jaques, Chase Kew, Jiaying Wu, Douglas Eck, Dale Schuurmans, and Aleksandra Faust. Joint attention for multi-agent coordination and social learning. *arXiv:2104.07750 [cs]*, August 2021. URL <http://arxiv.org/abs/2104.07750>. arXiv: 2104.07750.
- David Ha and Yujin Tang. Collective intelligence for deep learning: A survey of recent developments. *arXiv:2111.14377 [cs]*, Nov 2021. URL <http://arxiv.org/abs/2111.14377>. arXiv: 2111.14377.
- Yoshua Bengio. The consciousness prior. *arXiv:1709.08568 [cs, stat]*, December 2019. URL <http://arxiv.org/abs/1709.08568>. arXiv: 1709.08568.
- Michael S. A. Graziano. The attention schema theory: A foundation for engineering artificial consciousness. *Frontiers in Robotics and AI*, 4:60, November 2017. doi:10.3389/frobt.2017.00060. URL <http://journal.frontiersin.org/article/10.3389/frobt.2017.00060/full>.
- Bernard J. Baars. *A Cognitive Theory of Consciousness*. Cambridge University Press, July 1993. ISBN 9780521427432. Google-Books-ID: 7w6IYeJRqyoC.
- Stanislas Dehaene, Michel Kerszberg, and Jean-Pierre Changeux. A neuronal model of a global workspace in effortful cognitive tasks. *Proceedings of the National Academy of Sciences*, 95(24):14529–14534, November 1998. doi:10.1073/pnas.95.24.14529. URL <https://www.pnas.org/content/95/24/14529>.
- Anirudh Goyal, Aniket Didolkar, Alex Lamb, Kartikeya Badola, Nan Rosemary Ke, Nasim Rahaman, Jonathan Binas, Charles Blundell, Michael Mozer, and Yoshua Bengio. Coordination among neural modules through a shared global workspace. *arXiv:2103.01197 [cs, stat]*, March 2021. URL <http://arxiv.org/abs/2103.01197>. arXiv: 2103.01197.

- Mingde Zhao, Zhen Liu, Sitao Luan, Shuyuan Zhang, Doina Precup, and Yoshua Bengio. A consciousness-inspired planning agent for model-based reinforcement learning. *arXiv:2106.02097 [cs]*, September 2021. URL <http://arxiv.org/abs/2106.02097>. arXiv: 2106.02097.
- Roger C. Conant and W. Ross Ashby. Every good regulator of a system must be a model of that system †. *International Journal of Systems Science*, 1(2):89–97, 1970. doi:10.1080/00207727008920220. URL <http://www.tandfonline.com/doi/abs/10.1080/00207727008920220>.
- Andrew I. Wilterson and Michael S. A. Graziano. The attention schema theory in a neural network agent: Controlling visuospatial attention using a descriptive model of attention. *Proceedings of the National Academy of Sciences*, 118(33), Aug 2021. doi:10.1073/pnas.2102421118. URL <https://www.pnas.org/content/118/33/e2102421118>.
- J. P. Mitchell. Activity in right temporo-parietal junction is not selective for theory-of-mind. *Cerebral Cortex*, 18(2):262–271, February 2008. doi:10.1093/cercor/bhm051. URL <https://academic.oup.com/cercor/article-lookup/doi/10.1093/cercor/bhm051>.
- Michael S. A. Graziano and Sabine Kastner. Human consciousness and its relationship to social neuroscience: A novel hypothesis. *Cognitive Neuroscience*, 2(2):98–113, 2011. doi:10.1080/17588928.2011.565121. URL <http://www.tandfonline.com/doi/abs/10.1080/17588928.2011.565121>.
- Rodney A. Brooks. Intelligence without representation. *Artificial Intelligence*, 47(1):139–159, Jan 1991. doi:10.1016/0004-3702(91)90053-M. URL <https://www.sciencedirect.com/science/article/pii/000437029190053M>.
- Allen Newell and Herbert A. Simon. Computer science as empirical inquiry: symbols and search. *Communications of the ACM*, 19(3):113–126, March 1976. doi:10/dnhcsn. URL <https://doi.org/10.1145/360018.360022>.
- Tim van Gelder. The dynamical hypothesis in cognitive science. *Behavioral and Brain Sciences*, 21(5):615–628, 1998. doi:10/dk8nfn. URL https://www.cambridge.org/core/product/identifier/S0140525X98001733/type/journal_article.
- Tim Smithers. *Are Autonomous Agents Information Processing Systems?*, page 123–162. Routledge, 1 edition, May 2018. ISBN 9781351001885. doi:10.4324/9781351001885-4. URL <https://www.taylorfrancis.com/books/9781351001878/chapters/10.4324/9781351001885-4>.
- Xabier E. Barandiaran. Autonomy and enactivism: Towards a theory of sensorimotor autonomous agency. *Topoi*, 36(3):409–430, 2017. doi:10/ghh82k. URL <http://link.springer.com/10.1007/s11245-016-9365-4>.
- Karl Friston, James Kilner, and Lee Harrison. A free energy principle for the brain. *Journal of Physiology-Paris*, 100(1):70–87, July 2006. doi:10/dxgwt3. URL <https://www.sciencedirect.com/science/article/pii/S092842570600060X>.
- Maxwell JD Ramstead, Michael D Kirchhoff, and Karl J Friston. A tale of two densities: active inference is enactive inference. *Adaptive Behavior*, 28(4):225–239, August 2020. doi:10/gf97c4. URL <https://doi.org/10.1177/1059712319862774>.
- Andy Clark. Whatever next? predictive brains, situated agents, and the future of cognitive science. *Behavioral and Brain Sciences*, 36(3):181–204, June 2013. doi:10/f4xkv5. URL <https://www.cambridge.org/core/journals/behavioral-and-brain-sciences/article/whatever-next-predictive-brains-situated-agents-and-the-33542C736E17E3D1D44E8D03BE5F4CD9>.
- Jakob Hohwy. *The Predictive Mind*. Oxford University Press, November 2013. ISBN 9780199682737. Google-Books-ID: z7gVDAAAQBAJ.
- Linda Smith and Michael Gasser. The development of embodied cognition: six lessons from babies. *Artificial Life*, 11(1–2):13–29, 2005. doi:10.1162/1064546053278973.
- Evan Thompson and Francisco J. Varela. Radical embodiment: neural dynamics and consciousness. *Trends in Cognitive Sciences*, 5(10):418–425, October 2001. doi:10/c75qx5.
- Lawrence W. Barsalou, Paula M. Niedenthal, Aron K. Barbey, and Jennifer A. Ruppert. Social embodiment. *Psychology of Learning and Motivation - Advances in Research and Theory*, page 43–92, 2003. doi:10.1016/S0079-7421(03)01011-9. URL <https://experts.illinois.edu/en/publications/social-embodiment>.
- Eric Deng, Bilge Mutlu, and Maja Mataric. Embodiment in socially interactive robots. *Foundations and Trends in Robotics*, 7(4):251–356, 2019. doi:10.1561/23000000056. URL <http://arxiv.org/abs/1912.00312>. arXiv: 1912.00312.
- Mengsen Zhang, Guillaume Dumas, JA Scott Kelso, and Emmanuelle Tognoli. Enhanced emotional responses during social coordination with a virtual partner. *International Journal of Psychophysiology*, 104:33–43, 2016.
- Hatice Kose-Bagci, Ester Ferrari, Kerstin Dautenhahn, Dag Sverre Syrdal, and Chrystopher L. Nehaniv. Effects of embodiment and gestures on social interaction in drumming games with a humanoid robot. *Advanced Robotics*, 23(14):1951–1996, January 2009. doi:10.1163/016918609X12518783330360. URL <https://doi.org/10.1163/016918609X12518783330360>.
- Agrim Gupta, Silvio Savarese, Surya Ganguli, and Li Fei-Fei. Embodied intelligence via learning and evolution. *Nature Communications*, 12(1):5721, October 2021.

- doi:10.1038/s41467-021-25874-z. URL <https://www.nature.com/articles/s41467-021-25874-z>.
- Kwan Lee. Are physically embodied social agents better than disembodied social agents?: The effects of physical embodiment, tactile interaction, and people's loneliness in human–robot interaction. *International Journal of Human-Computer Studies*, 64(10):962–973, October 2006. doi:10.1016/j.ijhcs.2006.05.002. URL <https://www.sciencedirect.com/science/article/abs/pii/S1071581906000784>.
- Guillaume Dumas, Julien Laroche, and Alexandre Lehmann. Your body, my body, our coupling moves our bodies. *Frontiers in human neuroscience*, 8:1004, 2014a.
- Michael Dennis, Natasha Jaques, Eugene Vinitsky, Alexandre Bayen, Stuart Russell, Andrew Critch, and Sergey Levine. Emergent complexity and zero-shot transfer via unsupervised environment design. *arXiv:2012.02096 [cs]*, Feb 2021. URL <http://arxiv.org/abs/2012.02096>. arXiv: 2012.02096.
- AM Turing. Computing machinery and intelligence. *Mind*, 59(236):433–460, 1950.
- Guillaume Dumas, Gonzalo C de Guzman, Emmanuelle Tognoli, and JA Scott Kelso. The human dynamic clamp as a paradigm for social interaction. *Proceedings of the National Academy of Sciences*, 111(35):E3726–E3734, 2014b.
- Alan Mathison Turing et al. On computable numbers, with an application to the entscheidungsproblem. *J. of Math*, 58(345-363):5, 1936.
- AM Turing. The chemical basis of morphogenesis. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, 237(641):37–72, 1952.