

A Theory of Social Agency for Human-Robot Interaction

Ryan Blake Jackson^{1,*} and Tom Williams¹

¹MIRRORLab, Department of Computer Science, Colorado School of Mines, Golden, CO, USA

Correspondence*:
Ryan Blake Jackson
rbjackso@mines.edu

2 ABSTRACT

3 Motivated by inconsistent, underspecified, or otherwise problematic theories and usages of
4 *social agency* in the HRI literature, and leveraging philosophical work on *moral agency*, we
5 present a theory of social agency wherein a social agent (a thing with social agency) is any *agent*
6 capable of *social action* at some *level of abstraction*. Like previous theorists, we conceptualize
7 *agency* as determined by the criteria of interactivity, autonomy, and adaptability. We use the
8 concept of *face* from politeness theory to define *social action* as any action that threatens or
9 affirms the face of a *social patient*. With these definitions in mind, we specify and examine the
10 levels of abstraction most relevant to HRI research, compare notions of social agency and the
11 surrounding concepts at each, and suggest new conventions for discussing social agency in our
12 field.

13 **Keywords:** social agency, human-robot interaction, moral agency, politeness theory, levels of abstraction

1 INTRODUCTION AND MOTIVATION

14 The terms “social agency” and “social agent” appear commonly within the human-robot interaction
15 (HRI) research community. From 2011 to 2020, these terms appeared in at least 45 papers at ACM/IEEE
16 International Conference on HRI alone¹, with more instances in related conferences and journals. Given
17 the frequency with which these terms are used in the HRI community, one might expect the field to have
18 established agreed upon definitions to ensure precise communication. However, when these terms are used,
19 they are often not explicitly defined, and their use frequently varies in important but subtle ways, as we will
20 discuss below. Most HRI research is not concerned with exploring the entire philosophy of agency to find
21 a theory that fits their study. As we show in Section 1.3, it is therefore common to simply use terms like
22 “social agency” without espousing a particular concrete definition and move on under the assumption that it
23 is clear enough to the reader what is meant. This may be fine within any individual paper, but confusion
24 arises when different papers in the same research area use the same term with different meanings. We
25 seek to formalize social agency in accordance with the existing underspecified usage because (1) having
26 a rigorously specified definition for the term will help create common ground between researchers, help
27 new researchers understand the vernacular of the community, and provide writing guidelines for HRI
28 publications concerning social agency; and (2) attempting to redefine social agency in a substantially
29 different way from existing habits of use would greatly hamper popular acceptance of the new definition.

¹ <https://dl.acm.org/action/doSearch?AllField=%22social+agent%22+%22social+agency%22&ConceptID=119235>

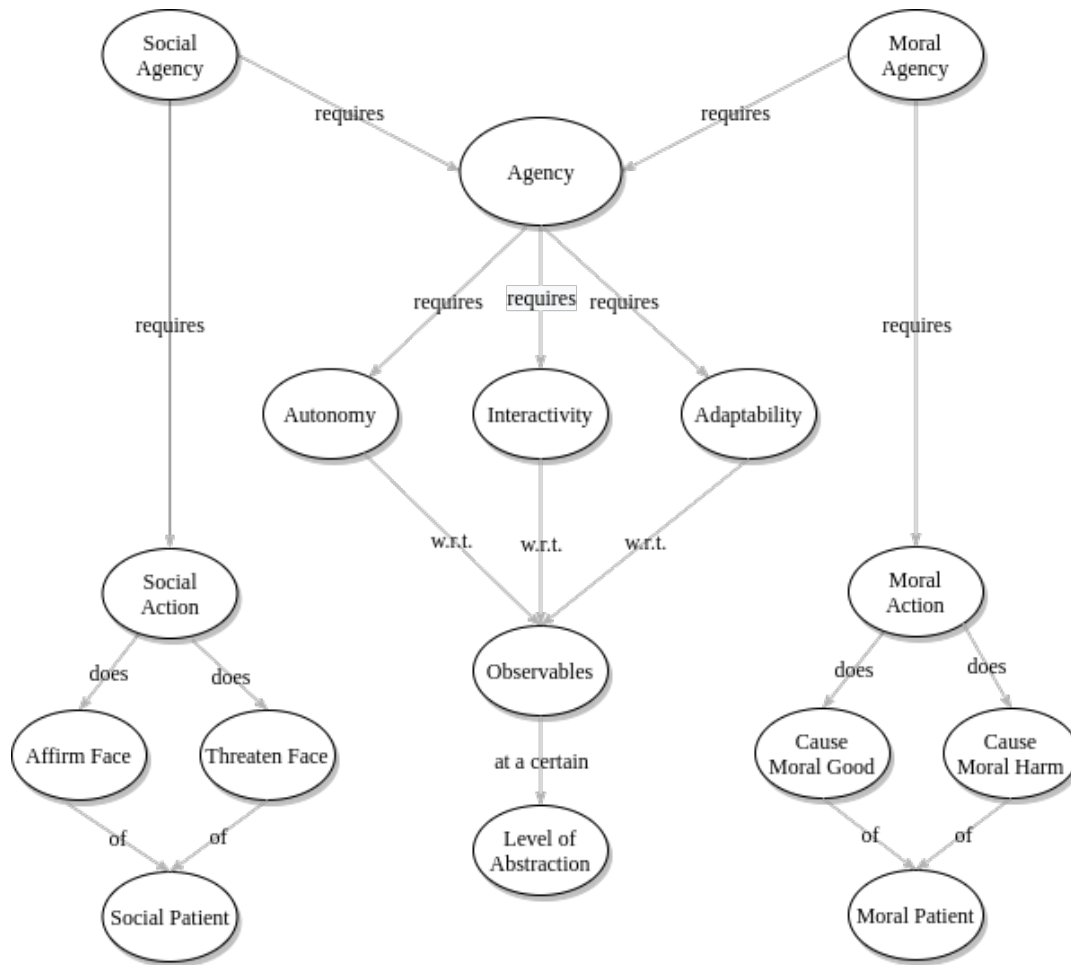


Figure 1. Concept Diagram visualizing the theory of Social Agency presented in this paper, and the core concepts combined to construct this theory.

30 We present a theory of social agency for HRI research (as visualized in Fig. 1) that deliberately aligns
 31 with and builds on other philosophical theories of robot agency. Specifically, we leverage insights from
 32 philosophers seeking to define *moral* agency in HRI. Moral agency provides an excellent analog to facilitate
 33 our discussion of social agency because it is an intimately related concept for which scholars have already
 34 developed rigorous definitions applicable to HRI, in a way that has not yet been done for social agency.

35 To design and justify our theory of social agency, we will first briefly survey existing definitions of social
 36 agency outside of HRI, and explain why those definitions are not well-suited for HRI. We will then survey
 37 theories of social agency from within HRI, and explain why those definitions are both inconsistent with
 38 one another and insufficient to cover the existing casual yet shared notion of social agency within our field.
 39 To illustrate this existing notion, we will then present a representative sample of HRI research that refers to
 40 social agency (without focusing on developing a definition thereof) to demonstrate how the greater HRI
 41 community's casual use of social agency differs from the more rigorous definitions and theories found
 42 within and beyond the field of HRI.

43 1.1 Social Agency Outside HRI

44 There are many different definitions of social agency from various disciplines including Psychology,
 45 Education, Philosophy, Anthropology, and Sociology. Providing an exhaustive list of these differing
 46 definitions is infeasible, but this section briefly summarizes a few representative definitions from different

47 fields to show that they are not well-suited to HRI and to illustrate the broader academic context for our
48 discussion of social agency.

49 Educational psychologists have used the term “social agency theory” to describe the idea that
50 computerized multimedia learning environments “can be designed to encourage learners to operate under
51 the assumption that their relationship with the computer is a social one, in which the conventions of
52 human-to-human communication apply” (Atkinson et al., 2005). Essentially, social agency theory posits
53 that the use of verbal and visual cues, like a more humanlike than overtly artificial voice, in computer-
54 generated messages can encourage learners to consider their interaction with the computer to be similar
55 to what they would expect from a human-human conversation. Causing learner attributions of social
56 agency is hypothesized to bring desirable effects, including that learners will try harder to understand the
57 presented material (Atkinson et al., 2005). In contrast, typically in HRI to be a social agent is humanlike in
58 that humans are social agents, but more human-likeness, particularly in morphology or voice, does not
59 necessarily imply more social agency. This theory also seems fundamentally concerned with social agency
60 creating a social partnership to facilitate learning, but we also view non-cooperative social behaviors, like
61 competition or argument, as socially agentic (Castelfranchi, 1998).

62 Other education researchers use the term social agency differently. For example, though Billett (2008)
63 does not explicitly define social agency (a practice that we will see is common in HRI literature as well),
64 they seem to view social agency as the capacity for the greater social world to influence individuals. This
65 concept contrasts with personal agency, which Billett defines explicitly as an individual’s intentional actions.
66 Personal and social agencies exert interdependent forces on the human worker as they negotiate their
67 professional development and lives. This notion of social agency that precludes it from being a property
68 held by a single individual, which does not seem to be how we use the term in HRI.

69 Scholars in education and social justice have also defined social agency as the extent to which individuals
70 believe that being active socio-politically to improve society is important to their lives, and the extent to
71 which individuals believe that they can / ought to alter power relations and structural barriers (Garibay,
72 2015, 2018). This definition is largely centered around value placed on prosocial behavior. In contrast, in
73 HRI we often apply the concept of social agency regardless of whether a robot is having any nontrivial
74 impact on society or is trying to do so. We also ascribe social agency regardless of what a robot believes or
75 values, or whether it can even believe or value anything.

76 Much of the discussion around agency in Anglo-American philosophy has revolved around intentionality,
77 but some influential anthropologists have centered not only intentionality in defining agency, but also the
78 power, motivation, and requisite knowledge to take consequential action (Gardner, 2016). Social agency,
79 then, could be understood as agency situated within a social environment, wherein agents produce and
80 reproduce the structures of social life, while also being influenced by those structures (and other material
81 conditions), particularly through the rules, norms, and resources that they furnish. Social agency here is
82 concerned with structures and relationships of power between actors. Other scholars in anthropology and
83 related fields have criticized this notion of agency, for, among other reasons, over-emphasizing the power
84 of the individual and containing values particular to men in the modern West. Some scholars that have
85 de-emphasized power and capacity have stated that intentions alone are what characterize an agent and
86 choices are the outcomes of these intentions, without necessarily qualitatively redefining the relationship
87 between agency and social agency (Gardner, 2016). These definitions, and other similar ones, are also
88 common in sociology and other social sciences. For reasons that we will argue below, we avoid “internal”
89 factors like intentionality, motivation, and knowledge in defining social agency for HRI. We are also not

90 concerned with whether robots have the power to act with broad social consequences since that does not
91 seem important to HRI researcher's usage of the term.

92 Anthropologists and archaeologists apply "social agency theory" to the study of artifactual tools and
93 technologies to understand the collective choices that were made during the manufacture and use of such
94 artifacts, the intentions behind those choices, the sociocultural underpinnings of those intentions, and the
95 effects that the technologies had on social structures and relations. In doing so, they commonly refer to the
96 social agency of technology or of technological practice to discuss the relationships between a technology
97 and the social structures and decisions of its manufacturers and users. For example, the choice to use
98 inferior local materials for tools rather than sourcing better materials through commerce given the material
99 means to do so can indicate constraining social structures outweighing the enabling economic structures
100 (Dobres and Hoffman, 1994; Gardner, 2016). Contrastingly, in HRI robots are discussed as having social
101 agency in and of themselves, separate from that of the humans that make and use them. Social robots are
102 also attributed social agency without really being embedded in the same broader social structures as their
103 human interactants, though it is likely that they will be increasingly as the field progresses.

104 Scholars in Sociology have also conceptualized agency as the constructed authority, responsibility, and
105 legitimated capacity to act in accordance with abstract moral and natural principles. Modern actors (e.g.,
106 individuals, organizations, and national states) have several different sorts of agency. Agency for the
107 self involves the tendency of an actor towards elaborating its own capacities in accordance with wider
108 rationalized rules that define its agency, even though such efforts are often very far removed from its
109 immediate raw interests. For example, organizations often develop improved information systems toward
110 no immediate goal. Agency for other actors involves opining, collaborating, advising, or modeling in
111 service of others. Agency for nonactor entities is the mobilization for culturally imagined interests of
112 entities like ecosystems or species. Finally, agency for cultural authority describes how, in exercising
113 any type of agency, the actor assumes responsibility to act in accordance with the imagined natural and
114 moral law. At the extreme, actors can represent pure principle rather than any recognized entity or interest.
115 However, for the modern actor, being an agent is held in dichotomy with being a principal, where the
116 principal "has goals to pursue or interests to protect, [and] the agent is charged to manage this interestedness
117 effectively, but in tune with general principles and truths." In other words, the principal is concerned with
118 immediate raw interests, while the agent is concerned with higher ideals. For example, the goals of a
119 university as principal are to produce education and research at low cost, whereas the goals of the university
120 as agent include having the maximum number of brilliant (expensive) professors and the maximum number
121 of prestigious programs. The same tension manifests in individuals as classic psychological dualisms (e.g.,
122 short-term vs. long-term interests) By this duality, highly agentic features like opinions and attitudes can
123 be decoupled from behaviors, actions, and decisions (Meyer and Jepperson, 2000).

124 Social agency, within this body of work, refers to the social standardization and scriptedness of agency,
125 and to how agency dynamics permeate and shape social structure. In a society of social agents, each
126 individual or organization acts in accordance with their socially prescribed and defined agency, which is
127 akin to the ideals defining their social role. In general terms, "the actorhood of individuals, organizations,
128 and national states [is] an elaborate system of social agency..." wherein actors routinely shift between
129 agency for the self and otherhood for the generalized agency of the social system. Individuals share in
130 the general social agency of the system, negotiating the bases for their own existence via the rules and
131 definitions of the broader system. This general social agency can function as the capacity for collective
132 agentic action (Meyer and Jepperson, 2000). This understanding of agency as an upholding of higher
133 ideals, principles, and truths (and social agency as the collective version of this), often in conflict with

134 baser self-interested principalhood, is so different from conceptions of agency and social agency in HRI as
135 to be essentially completely disjoint concepts. As we will illustrate below, agency in HRI is not (to our
136 knowledge) discussed in duality with the notion of a principal, and social agency is not understood as a
137 collective version of individual agency.

138 In presenting the definitions in this section, we do not intend to suggest that other fields have reached
139 some sort of internal consensus regarding social agency or perfect consistency in its usage. Like in HRI,
140 there appears to be ongoing conversation and sometimes disagreement about social agency within many
141 fields, though the HRI-specific branch of this conversation seems relatively nascent. For example, there are
142 ongoing debates in anthropology about whether (social) agency is an essential property of individuals, or
143 somehow exists only in the relationships between individuals. Likewise, there are differing opinions within
144 and between social science research communities about whether nonhuman entities can have (social) agency
145 Gardner (2016). Unfortunately, we cannot present all perspectives here, nor can we really present the full
146 detail and nuance of some of the perspectives that we *have* presented. What we hope to have indicated is
147 that definitions of social agency from other fields, though academically rigorous and undoubtedly useful
148 within their respective domains, are, for various reasons, neither intended nor suitable for the unique role
149 of social agency in HRI, and an HRI-specific definition is needed.

150 1.2 Theories of Social Agency in HRI

151 A number of theories of Social Agency have been defined within the HRI community to address the
152 unique perspective of our field. Many of these grew out of foundational work on Social Actors from Nass
153 et al. (1994), which suggested that humans naturally perceive computers with certain characteristics (e.g.,
154 linguistic output) as social *actors*, despite knowing that computers do not possess feelings, “selves”, or
155 human motivations (Nass et al., 1994). This perception leads people to behave socially towards machines
156 by, for example, applying social rules like politeness norms to them (Nass et al., 1994; Jackson et al.,
157 2019). It is perhaps unsurprising that this human propensity to interact with and perceive computers
158 in fundamentally social ways extends strongly to robots, which are often deliberately designed to be
159 prosocial and anthropomorphised. While Nass et al.’s work establishing the theory that humans naturally
160 view computers as social *actors* did not call computers “social agents” or refer to the “social agency” of
161 computers, it nevertheless established that the human-computer relationship is fundamentally social, and
162 laid the groundwork for much of the discussion of sociality and social agency in HRI today. In this section
163 we will discuss four rigorously defined theories of Social Agency in HRI.

164

165 Nagao and Takeuchi

166 At around the same time that Nass and colleagues introduced their “Computers As Social Actors” (CASA)
167 paradigm (Nass et al., 1994), Nagao and Takeuchi (1994) made one of the earliest references to computers
168 as *social agents*. In describing their approach to social interaction between humans and computers, Nagao
169 and Takeuchi argue that a computer is a social agent if it is both social and autonomous. These authors
170 define socialness as multimodal communicative behavior between multiple individuals. Nagao and Takeuchi
171 initially define autonomy as “[having] or [making] one’s own laws,” but later clarify that “an autonomous
172 system has the ability to control itself and make its own decisions.” We will see throughout this paper
173 that sociality and autonomy remain central to our discussion of social agency today, but not necessarily as
174 defined by these authors.

175 Nagao and Takeuchi also define a social agent as “any system that can do social interaction with humans,”
176 where a “social interaction” (1) involves more than two participants, (2) follows social rules like turn taking,
177 (3) is situated and multimodal, and (4) is active (which might be better understood as mixed initiative).

178 Some of these requirements, including at least the involvement of more than two participants and mixed
179 initiativity, seem unique to this theory. Nagao and Takeuchi also differentiate their “social interactions”
180 from problem solving interactions, though we believe, and see in the HRI literature, that task-oriented
181 interactions can be social and take place among social agents.

182 Pollini

183 Pollini (2009) presents a theory that is less concerned with modality of interaction or type of robot
184 embodiment, focusing instead on the role of human interactants in constructing a robot’s social agency. For
185 Pollini, robotic social agents are both physically and socially situated, with the ability to engage in complex,
186 dynamic, and contingent exchanges. Social agency, then, arises as the outcome of interaction with (human)
187 interlocutors, as “the ability to act and react in a goal-directed fashion, giving contingent feedback and
188 predicting the behavior of others.” We see the goal-directedness in this definition as loosely analogous to the
189 notion of autonomy that is centered in other theories. In contrast to those theories, however, Pollini considers
190 social agency as a dynamic and emergent phenomenon constructed collectively within a socially interacting
191 group of autonomous actors, rather than as an individual attribute separately and innately belonging to the
192 entities that comprise a social group. This presents a useful framing for understanding the social agency
193 of multi-agent organizations like groups and teams. However, this multi-agent perspective prevents this
194 definition from aligning with common references in HRI to the “social agency” of an individual robot.
195 Nonetheless, some degree of autonomous behavior, interaction, perception, and contingent reaction must
196 clearly remain central to our discussion of social agency.

197 Pollini also opines that “social agency is rooted in fantasy and imagination.” It seems that humans’
198 *attribution* of social agency may be tied to the development of imagination during childhood, leading
199 Pollini to argue that people can “create temporary social agents” of almost anything with which they have
200 significant contact, including toys like dolls, tools like axes, and places like the home. This leads them
201 to the question “what happens when such ‘entities-by-imagination’ also show autonomous behavior and
202 contingent reactions, and when they exist as social agents with their own initiative?” However, we argue
203 that axes, dolls, and places actually *cannot* be social agents, at least not in the way that the typical HRI
204 researcher means when they call a robot (or human) a social agent, since robots can conditionally take
205 interactional behavior, which we believe is necessary for social agency.

206 Finally, Pollini argues that agency-specific cues embedded in robots (e.g., contingent behavior) are
207 insufficient by themselves for creating social agency, and that social agency, rather, is negotiated between
208 machines and their human interactants via a process of interpretation, attribution, and signification. This
209 process involves interpreting a machine’s behavior as meaningful and explicative, and then attributing
210 social agency based on the signification of that behavior as meaningful, which may also involve attributing
211 internal forces like intentions and motivations. This means that, through this process, things with simple
212 behaviors like cars or moving shapes on a screen can end up being ascribed social agency. Again, however,
213 we see a fundamental difference between these examples and social robots, which can actually deliberately
214 manifest meaningful and explicative behaviors. We interpret this discussion as circling the distinction
215 between “actual” and “perceived” social agency that we will discuss below.

216

217 Levin, Adams, Saylor, and Biswas

218 Though much of the HRI literature exploring the standalone concept of *agency* is beyond the scope of this
219 work as it focuses on the agency of machines without centering notions of *sociality*, the theory of agency
220 from Levin et al. (2013) is relevant here because it explores attributions of agency specifically during social
221 human-robot interactions. Levin et al. argue that people’s first impulse is to strongly differentiate the agency

222 of humans and nonhumans, and that people only begin to equate the two with additional consideration
223 (e.g., when prompted to do so by the robot defying initial expectations). They also describe how simple
224 robot behavioral cues like the naturalness of movement or gaze can influence people's attribution of agency
225 to robots, as well as states and traits of the human attributor, like loneliness. Like some previous theories,
226 Levin et al. center goal-orientedness and intentionality in their account of agency. However, they include
227 not only behavioral intentionality, which we saw in other theories (Pollini, 2009), but also intentionality in
228 cognition. Their example of this cognitive intentionality is drawing ontological distinctions between types
229 of objects based on their use rather than their perceptual features.

230 Alač

231 Finally, Alač (2016) presents a theory in which multimodal interaction, situatedness, and materiality are
232 important to a robot's social agency, and justifies this theory with an observational study of a robot in a
233 classroom. Alač frames robot agenthood as coexisting with the contrasting status of "thing," with agentic
234 features entangled in an interplay with a robot's thing-like materiality. However, Alač moves away from
235 discussing a robot's social nature as an intrinsic and categorical property that resides exclusively in the
236 robot's physical body or programming, instead seeing robot sociality as enacted and emergent from how
237 a robot is experienced and articulated in interactions. To Alač, the socially agentic facets of a robot are
238 evident in the way it is treated by humans, focusing on proxemic and haptic interaction patterns and
239 linguistic framing (e.g., gendering the robot) in group settings. Our work can augment ethnography-based
240 theories like this one by exploring (1) the features of the *robot's* behavior that give rise to perceptions of
241 social agency, (2) what concepts constitute such perceptions, and (3) exactly what such perceptions imply.
242 In other words, we focus on what social agency *is*, rather than on human behaviors that indicate ascription
243 thereof.

244

245 1.3 Notions of Social Agency in HRI

246 While in the previous section we discussed rigorously defined theories of social agency, much of the HRI
247 literature that engages with social agency does not actually connect with those theories. In this section,
248 we will thus explore the ways in which HRI researchers casually refer to social agency without focusing
249 on developing or defining a formal theoretical account of it. Our goals in doing so are to (1) illustrate
250 that notions of social agents and agency are commonly applied within the HRI research community, (2)
251 provide examples of *how* these terms are used, and demonstrate important qualitative differences among
252 the entities to which these terms are applied, (3) show that the existing theories defined in the previous
253 section do not capture the common parlance usage of "social agency" among HRI researchers, and (4) lay
254 the groundwork for developing a theory that does accommodate these usages.

255 There are many papers that refer to robots as social agents without mentioning or dealing with *social*
256 *agency* per se. The term social agent is widely applied to entities that are both embodied (Heerink et al.,
257 2010; Lee et al., 2012; Luria et al., 2016; Westlund et al., 2016) and disembodied (Lee et al., 2006; Heerink
258 et al., 2010); remote controlled by humans (Heerink et al., 2010; Lee et al., 2012; Westlund et al., 2016)
259 and self-controlled (Heerink et al., 2010); task-oriented (Heerink et al., 2010; Lee et al., 2012) and purely
260 social (Lee et al., 2006); anthropomorphic (Heerink et al., 2010; Lee et al., 2012), zoomorphic (Lee et al.,
261 2006; Heerink et al., 2010; Westlund et al., 2016), and mechanomorphic (Heerink et al., 2010; Luria et al.,
262 2016); mobile (Heerink et al., 2010; Lee et al., 2012) and immobile (Heerink et al., 2010; Luria et al.,
263 2016); and able to communicate with language (Heerink et al., 2010; Lee et al., 2012) and unable to do
264 so (Lee et al., 2006; Luria et al., 2016). Any theory of social agency for HRI, then, should either encompass
265 this diversity of social agents or account for ostensible misattributions of social agency. However, the

266 theories we have examined, which emphasize embodiment (Nagao and Takeuchi, 1994; Alač, 2016),
267 language (Nagao and Takeuchi, 1994), and self-control or intentionality (Pollini, 2009; Levin et al., 2013),
268 exclude usages that are apparently common in HRI research.

269 Of course, one could argue that casual references to robots as “social agents” are synonymous to
270 references to robots as “social actors,” and that such references do not actually have anything to do with
271 the agentic nature of the robot. By this argument, the existing theoretical work on social agency in HRI
272 would best be understood as investigating a completely separate topic from social agents. This reasoning,
273 however, would result in a confusing state-of-affairs in which social agency is not a prerequisite for being
274 a social agent, with the two topics unrelated except by the general connection to social interaction. We
275 therefore assume that a social agent must be a thing with social agency, and that these two terms must be
276 tightly and logically related. A clear conception of social agency is thus a prerequisite for the study of
277 social agents. However, much of the work in HRI that concerns social agency does not focus on rigorously
278 defining it. Indeed, some of these studies do not explicitly provide their definition of social agency at all.

279 An illustrative example of a casually referenced “social agent” is the “Snackbot” developed by Lee et al.
280 (2012). The anthropomorphic Snackbot had real interactions with many humans over the course of multiple
281 months as a snack delivery robot. The robot’s movement was self-controlled, but a human teleoperator
282 hand-selected its delivery destinations. The human operator also remotely controlled the robot’s head and
283 mouth movements and the robot’s speech, by selecting from a number of pre-made scripts, both purely
284 social and task-oriented. We will refer back to this example in Section 2.

285 In their investigation of how cheating affects perceptions of social agency, Ullman et al. (2014) used
286 perceptions of trustworthiness, intelligence, and intentionality as indicators of perceptions of social agency
287 in an anthropomorphic robot. Using intentionality as a proxy for social agency aligns directly with several
288 of the theories that we described in Section 1.2 (Pollini, 2009; Levin et al., 2013). Intelligence and
289 trustworthiness, however, seem less closely related to social agency, and trustworthiness is explicitly not an
290 aspect of social agency in theories that discuss competition and uncooperative behavior as inherently social
291 actions (Castelfranchi, 1998).

292 Baxter et al. (2014) also study attributions of social agency to robots without explicitly defining the term,
293 and measure it via a different proxy: human gaze behavior. This proxy does not obviously align with any of
294 the theories of social agency discussed above. Although it is possible that gaze could be a good proxy for
295 some definition of social agency (or the ascription thereof), further empirical work would be needed to
296 establish that relationship.

297 Straub (2016) adopt yet another definition of social agency in their investigation of the effects of social
298 presence and interaction on social agency ascription. In their study, social agents are characterized as
299 “having an ‘excentric positionality,’ equipped with (a) an ability to distinguish themselves, their perceptions
300 as well as their actions from environmental conditions (embodied agency), (b) the ability to determine
301 their actions and perceptions as self-generated, (c) having the ability to define and relate to other agents
302 equipped with the same features of (a) and (b), along with (d) defining their relationship to other agents
303 through reciprocal expectations toward each other (‘excentric positioned’ alter ego).”

304 This definition, particularly part b, is somewhat ambiguous. One interpretation is that the robot simply
305 needs to distinguish its own actions from the actions of others, and know that it is the cause for the effects
306 of its actions; if the robot moves its arm into a cup, then it is the source for both the movement of the arm
307 and the movement of the cup. However, this seems more like the robot knowing that its actions’ effects
308 are self-generated and that it was the one that acted, rather than viewing the choice to act or the genesis

309 of the action itself as self-generated. Another interpretation, which is similar to some of the definitions
310 of social agency discussed in Section 1.1, is that seeing an action as self-generated requires the robot to
311 understand its choice to act, perceive that choice as its own, and believe that it could have acted differently.
312 This definition appears to require some form of consciousness or experience of free will, and is thus not
313 well-suited to HRI. Straub uses human behavioral proxies, like eye contact, mimicry, smiles, and utterances,
314 to measure ascriptions of social agency to robots (with more of these behaviors indicating more ascribed
315 social agency), but such behavioral proxies do not measure all components of their definition.

316 Ghazali et al. (2019) study the effects of certain social cues (emotional intonation of voice, facial
317 expression, and head movement) on ascriptions of social agency. Professedly inspired by research in
318 educational psychology described above (Atkinson et al., 2005), they define social agency as “the degree
319 to which a social agent is perceived as being capable of social behavior that resembles human-human
320 interaction,” and then measure it by collecting participant assessments of the extent to which the robot was
321 “real” and “like a living creature.” Roubroeks et al. (2011) use the exact same definition of social agency as
322 Ghazali et al. (2019) in their investigation of psychological reactance to robots’ advice or requests, but
323 operationalize it differently. Although they did not attempt to measure social agency, they did seek to
324 manipulate it by varying robot presentation, presenting a robot’s advice as either text alone, text next to a
325 picture of the robot, or a video of the robot saying the advice.

326 This definition seems problematically circular in that it defines social agency by the degree to which
327 a social agent does something, without defining what it means to be a social agent. We also argue that
328 Ghazali et al.’s chosen measures do not clearly align with the formal definitions of social agency proposed
329 above, nor with Ghazali et al.’s stated definition. Moreover, this conceptualization excludes a large number
330 of robots that the HRI literature calls social agents, and focuses on factors that many theories de-emphasize
331 (e.g., livingness and human likeness). This example in particular shows that disparate definitions of social
332 agency currently exist in the HRI literature, leading to confusion when authors underspecify or neglect to
333 specify a definition.

334 Other work from Ghazali et al. (2018) on the relationship between social cues and psychological reactance
335 centers the concepts of “social agent” and “social agency” explicitly, using the terms over 100 times in
336 reference to robots and computers. However, the authors do not expressly provide any definition for those
337 terms, despite ostensibly manipulating social agency in an experiment. Implicitly, the authors appear to
338 follow their definition described above, with more humanlike superficial behavior (e.g., head/eye movement
339 and emotional voice intonation) being considered more socially agentic, while the semantic content and
340 illocutionary force of all utterances was kept constant across social agency conditions. However, Ghazali
341 et al. (2018) also seem to consider the capacity to threaten others’ autonomy as a critical feature of social
342 agency, since they measure perceived threat to autonomy as a manipulation check on social agency (though
343 the social agency manipulation did not significantly impact perceived threat to autonomy). This choice
344 was not extensively justified. As discussed in Section 2.2, perceived threat to autonomy is strongly related
345 to (negative) face threat, which we view as important to social agency. However, as we will discuss, the
346 capacity to threaten face is far broader than the capacity to threaten autonomy as measured by Ghazali et al.
347 (2018).

348 To summarize, we have discussed several conflicting theories and usages of social agency in HRI, which,
349 to varying extents: (a) exclude common uses of the term “social agency” by being too restrictive, (b)
350 include objects that nearly all researchers would agree are neither social nor agentic, (c) focus on factors
351 that do not seem relevant to social agency in most pertinent HRI work, or (d) conflate other concepts (like
352 livingness or human-likeness) with social agency as it seems commonly understood. In addition, we have

Table 1. Summary of terms that are important to our concept of social agency.

Term	Definition
Level of Abstraction (LoA)	A collection of observables describing an entity (Floridi and Sanders, 2004; Floridi, 2008). A user's LoA for a robot includes movement, speech, morphology, etc., while the developer's LoA also includes the algorithms controlling the robot.
Agent	Anything possessing the three criteria of interactivity, autonomy, and adaptability.
Interactivity	The capacity to act on the environment and to be acted upon by the environment (Floridi and Sanders, 2004).
Autonomy	The capacity to change state without direct response to interaction (Floridi and Sanders, 2004).
Adaptability	The capacity for interaction to change the system's state transition rules. The capacity to "learn" from interaction (Floridi and Sanders, 2004).
Social Agent	Anything capable of taking social action at the LoA under consideration.
Social Action	Any act that threatens or affirms an other's face. Analogous to moral action (doing harm/good to an other).
Social Patient	Anything that can be a recipient of social action, i.e., anything with face.
Face	The public self-concept (meaning self-concept existing in others) that all members of society want to preserve and enhance for themselves. Negative face: an individual's claim to freedom of action and freedom from imposition. Positive face: an individual's self-image and wants, and the desire that these be approved of by others (Brown and Levinson, 1987).

353 shown examples of the diversity of uses of the term "social agency" in the HRI research literature. We now
 354 contribute our own theory of social agency, with the specific intention of accommodating the HRI research
 355 community's existing notions of social agency.

2 A THEORY OF SOCIAL AGENCY FOR HRI

356 In this section, we propose a formal theory of social agency for HRI to address the challenges and
 357 limitations discussed in the previous sections. Our key arguments are: (1) social agency may be best
 358 understood through parallels to moral agency; (2) considering various levels of abstraction (LoAs) is
 359 critical for theorizing about any kind of agency; (3) a social agent can be understood as something with
 360 agency that is capable of social action; (4) social action is grounded in face; and (5) social and moral
 361 agency are related yet independent.

362 To best understand social agency, we draw parallels to recent work on moral agency. Not only are
 363 the concepts centered in theories of social agency discussed in Section 1.2 (e.g., autonomy, contingent
 364 behavior, and intentionality) also centered in many theories of moral agency, but the moral agency of
 365 robots and other artificial actors has also received a more rigorous treatment than social agency in the
 366 HRI literature. The moral agency literature thus represents a valuable resource for constructing a parallel
 367 theory of social agency. Furthermore, the two concepts of moral and social agency are inexorably linked,
 368 representing the two halves of interactional agency. They provide congruent relationships to (and means
 369 of understanding) moral/social norms and are key to our most foundational understandings of interaction.
 370 Given these similarities and connections, parallel understandings of the two concepts are not only intuitive

371 but necessary, and we see no reason to attempt to define moral and social agency completely separately.
372 For our purposes, we will leverage the moral agency theory of Floridi and Sanders (2004), but note that, as
373 with social agency, there is not yet consensus among scholars as to a single canonical definition of moral
374 agency, prompting ongoing debate (Johnson and Miller, 2008).

375 2.1 Agency and Levels of Abstraction

376 Because of historical difficulties in defining necessary and sufficient conditions for agenthood that are
377 absolute and context-independent, Floridi and Sanders (2004) take analysis of *levels of abstraction* (LoAs)
378 (Floridi, 2008) as a precondition for analysis of agenthood. A LoA consists of a collection of observables,
379 each with a well-defined set of possible values or outcomes. An entity may be described at a range of
380 LoAs. For a social robot, the observables defining an average user's LoA might only include the robot's
381 behavior and other external attributes, like robot morphology and voice. In contrast, the robot developer's
382 LoA would likely also include information internal to the robot, such as the mechanisms by which it
383 perceives the world, represents knowledge, and selects actions. Critically, a LoA must be specified before
384 certain properties of an entity, like agency, can be sensibly discussed, as a failure to specify a LoA invites
385 inconsistencies and disagreements stemming not from differing conceptions of agency but from unspoken
386 differences in LoA.

387 The "right" LoA for discussing and defining moral agency must accommodate the general consensus that
388 humans are moral agents. Floridi and Sanders (2004) propose a LoA with observables for the following
389 three criteria: interactivity (the agent and its environment can act upon each other), autonomy (the agent
390 can change its state without direct response to interaction), and adaptability (the agent's interactions can
391 change its state transition rules; the agent can "learn" from interaction, though this could be as simple as a
392 thermostat being set to a new temperature at a certain LoA). For the sake of simplicity, we will consider
393 LoAs consisting only of observations that a typical human could make over a relatively short temporal
394 window. These observables encompass some concepts that were important to the theories discussed in
395 Section 1.2 (e.g., autonomy and contingent behavior), and exclude others (e.g., teleological variables like
396 intentionality or goal-directedness), which we discuss more below. We also consider a criterion that was
397 *not* included in many theories for social agency, namely adaptability.

398 At the user's LoA, wherein the deterministic algorithms behind a robot's behavior are unobservable, the
399 robot is interactive, autonomous, and adaptable, and therefore is an agent. However, at the robot developer's
400 LoA (or what Floridi and Sanders (2004) call the "system LoA"), which includes an awareness of the
401 algorithms determining the robot's behavior, the robot loses the attribute of adaptability and is therefore
402 not an agent. These two LoAs will be important throughout the rest of this paper.

403 We argue that the distinction between these two LoAs (the user's and the developer's) explains why
404 some scholars have suggested conceptualizing and measuring "*perceived* moral agency" in machines as
405 distinct from moral agency itself. This notion of perceived moral agency would ostensibly capture "human
406 attribution of the status of a machine's agency and/or morality (independent of whether it actually has
407 agency or morality)" (Banks, 2019), and these authors could easily define "perceived social agency" the
408 same way.

409 Much of the impetus for defining these new concepts seems to be a desire to avoid the varied and
410 conflicting definitions for agency (and the social and moral variants thereof). Typically within HRI,
411 researchers are primarily concerned with how their robots are *perceived* by human interactants (the user's
412 LoA), and how those interactants might ascribe social agency to those robots. In that sense, perceived
413 social agency as a concept seems like a good way to allow researchers to focus on what they really care

414 about without getting mired in discussions of their robot’s “actual” agency, though it can still leave exactly
415 what is perceived as (socially) agentic underspecified.

416 However, as we saw in Section 1, authors seldom refer to perceived social agency (particularly since
417 we just defined it as parallel to perceived moral agency, which also does not seem to have caught on),
418 but rather use the unqualified term “social agency”. Thus, rather than attempting to enforce a change in
419 terminology, we propose that “perceived moral/social agency” should be understood as moral/social agency
420 at the robot user’s LoA, and “actual” moral/social agency is the corresponding notion at the developer’s
421 LoA. To illustrate, consider the SnackBot (Lee et al., 2012) described in Section 1.3. This robot was
422 largely remotely controlled by a human, but, at the snack orderer’s (user’s) LoA it is a social agent. At
423 the developer’s LoA, the robot is not an agent, but the system in aggregate might be considered socially
424 agentic since one of its constituent parts, the human, is a social agent in and of itself.

425 If SnackBot could manifest the same behavior without human input, it would still not be agentic at the
426 developer’s LoA insofar as its behavior is the direct result of deterministic algorithms that only act on its
427 state. However, it does intuitively *seem* more agentic, prompting us to consider another useful LoA: one
428 where we are aware of the general distributed system that controls a robot (in terms of software cognitive
429 architectural components, hardware components like cloud computing, and human teleoperators), but
430 not aware of the inner workings of each constituent part of that system. At this LoA, which we call the
431 “architecture LoA”, a robot that does its computation internally might be agentic, but a robot that is remote
432 controlled by either a person or another machine could not be an agent in and of itself. Hundreds of different
433 LoAs could be constructed with various degrees of detail regarding how a robot works, but this is largely
434 not constructive if humans are unlikely to ever view the robot from those LoAs. However, we believe that
435 the architecture LoA is realistic for many potential robot interactants, particularly those that might own
436 their own personal robots, or participants in laboratory HRI studies after the experimental debriefing.

437 At first glance, it would be easy to draw some parallels between our three main LoAs (developer’s,
438 architecture, and user’s) and Dennett’s three stances from which to view an entity’s behavior in terms of
439 mental properties (physical, design, and intentional) (Dennett, 1978). The user’s LoA in particular bears
440 loose resemblance to Dennett’s intentional stance because the user is aware only of the robot’s externally
441 observable behaviors, and may rationalize them by projecting internal states onto the robot. Likewise, our
442 architecture LoA is explicitly concerned with the parts comprising a robot’s distributed system and the
443 broad purpose of each constituent part, like the design stance, though it is not necessarily concerned with
444 the purpose of the robot itself as a whole. However, several key distinctions separate our three LoAs from
445 Dennett’s three stances. Most obviously, the developer’s LoA is unlike Dennett’s physical stance in that it is
446 concerned with the algorithms producing the robot’s behavior but not the specifics of their implementation
447 nor the hardware executing them.

448 More broadly, the three LoAs we have presented generally represent three of the *sets of information*
449 that real people are most likely to have regarding robots during HRI, but there is no reason for this set
450 of LoAs to be considered exhaustive, and no reason why our analysis of social agency cannot also apply
451 to any other LoA from which a person views a robot. In contrast, more rigidly tripartite approaches
452 to epistemological levelism, like Dennett’s, though readily formalized in terms of LoAs, contain an
453 implicit ontological commitment and corresponding presupposed epistemological commitment because
454 they privilege explanations over observable information (Floridi, 2008). That is not to say that such
455 approaches to multi-layered analysis are not interesting and illustrative to HRI. For example, many
456 researchers have explored whether humans naturally adopt the intentional stance towards robots and other
457 artificial entities like they do towards other humans (Marchesi et al., 2019; Perez-Osorio and Wykowska,

2019; Schellen and Wykowska, 2019; Thellman et al., 2017; Thellman and Ziemke, 2019). However, it seems intuitive that robot developers versus users might take the intentional stance towards robots to different extents and under different conditions, so we posit that a specification of LoA is helpful in considering Dennett's stances and other attitudinal stances in HRI in much the same way that it is to our discussion of social agency, rather than Dennett's stances being homeomorphic to the three LoAs most salient here.

Most current cognitive architectures are precluded from agency at the developer's LoA because any learning is typically a matter of updating the robot's state by the deterministic rules of its code, rather than an actual update to the rules for transitioning between states (Floridi and Sanders, 2004). This includes black-box systems, like deep neural networks, because their lack of interpretability comes from an inability to fully understand how the state results in behavior, not from actual adaptability. However, we accept that humans have adaptability, and see no theoretical reason why the same level of adaptability could not be implemented in future artificial agents. Of course, particularly within the theory of causal determinism, there exists an LoA wherein humans do not have agency if all human behavior is rooted in the physical and chemical reactions of molecules in the brain (a "physical" LoA *a la* Dennett). Regardless of the veracity of this deterministic point of view, it seems clear that no LoA precluding agency from existing in the universe as we know it is a useful LoA at which to discuss agency in HRI.

We adopt the above notion of LoA and criteria for agenthood from Floridi and Sanders (2004) for our theory of social agency for several reasons. First, different LoAs help us to account for different understandings of social agency in the HRI literature, as we saw in our discussion of "actual" versus "perceived" social agency. Second, we can explicitly avoid conflating moral/social *agency* with moral/social *responsibility* (i.e., worthiness of blame or praise), which is another discussion beyond the scope of this paper. Third, avoiding internal variables like intentionality, goal-directness, and free-will guarantees that our analysis is based only on what is observable and not on psychological speculation, since a typical robot user cannot observe these attributes in the internal code or cognitive processes of their robot; we thus prefer a phenomenological approach.

Having established an understanding of agency, we now need to define some notion of sociality congruent to Floridi and Sanders's notion of morality. However, we first want to point out that our justification for avoiding unobservable factors in defining and assessing (moral/social) agency parallels a similar argument from proponents of ethical behaviorism in defining and assessing the moral status of robots. Ethical behaviorism is an application of methodological behaviorism (as opposed to ontological behaviorism) to the ethical domain, which holds that a sufficient reason for believing that we have duties and responsibilities toward other entities (or that they have rights against us) can be found in their observable relations and reactions to their environment and ourselves. In other words, robots have significant moral status if they are roughly performatively equivalent to other entities that have significant moral status, and whatever is going on unobservably "on the inside" does not matter. This is not to say that unobservable qualia do not exist, nor do we deny that such qualia may be the ultimate metaphysical ground for moral status. However, the ability to ascertain the existence of these unobservable properties ultimately depends on some inference from a set of observable representations, so a behaviorist's point of view is necessary to respect our epistemic limits (Danaher, 2020). We agree with this reasoning. Our definition of social agency could be framed as a form of "social behaviorism" that specifies the behavioral patterns that epistemically ground social agency and, by considering LoAs, is sensitive to the behaviors that are actually observed, rather than the set of behaviors that are, in principle, observable.

501 Of course, avoiding attributes like intentionality or goal directedness in our definitions in favor of a
502 behaviorist approach does not completely free us from needing to rely on some form of inference. At
503 a minimum, making observations from sensory input requires the inference or faith that one's sensory
504 inputs correspond to some external reality. Likewise, our interactivity criterion for agency requires some
505 causal inference or counterfactual reasoning. For example, concluding that a robot can be acted on by
506 the environment requires the counterfactual inference that the robot's "response" to a stimulus would not
507 have occurred absent that stimulus. Unfortunately, requiring some inference is unavoidable. In light of
508 this, one could argue that it is equally reasonable and necessary to infer intention and goal directedness
509 from behavior. For example, pulling on a door handle might signal an intent to open the door with the
510 goal of getting into the building, even though the same behavior could also signal mindless programming
511 to tug on handles without representing goals or having intentions. We argue that the sensory and causal
512 inferences required by our framework are lesser epistemological leaps and more necessary and common
513 (and therefore more justifiable) than inferences about other agent's mental states like intentionality and
514 goals. We also emphasize that goals and intentions are apparently not important to social agency at the
515 developer's LoA, since we saw many robots referred to as social agents by their developers in Section 1.3
516 that did not internally represent goals or intentions, and their developers would have known that.

517 **2.2 Social Action Grounded in Face**

518 We now move on to developing a notion of sociality congruent to Floridi and Sanders's notion of morality.
519 For Floridi and Sanders (2004), any agent that can take moral action on another entity (e.g., do good or
520 evil; cause harm or benefit) is a moral agent. Any entity that can be the recipient of moral action (e.g.,
521 be harmed or benefited) is a moral patient. Most agents (e.g., people) are both moral agents and moral
522 patients, though research has indicated an inverse relationship between perceptions of moral agency and
523 moral patiency (e.g., neurodivergent adults are perceived more as moral patients and less as moral agents
524 than neurotypical adults) (Gray and Wegner, 2009).

525 Just as a *moral* agent is any agentic source of moral action, we can define a *social* agent as any agentic
526 source of social action. We ground our definition of social action in the politeness theoretic concept of
527 "face" (Brown and Levinson, 1987). Face, which consists of positive face and negative face, is the public
528 self-concept (meaning self-concept existing in others) that all members of society want to preserve and
529 enhance for themselves. Negative face is defined as an agent's claim to freedom of action and freedom from
530 imposition. Positive face consists of an agent's self-image and wants, and the desire that these be approved
531 of by others. A discourse act that damages or threatens either of these components of face for the addressee
532 or the speaker is a face threatening act. Alongside the level of imposition in the act itself, the degree of
533 face threat in a face threatening act depends on the disparity in power and the social distance between the
534 interactants. Various linguistic politeness strategies exist to decrease face threat when threatening face is
535 unavoidable or desirable. Conversely, a face affirming act is one that reinforces or bolsters face for the
536 addressee or speaker (though our focus will be on the addressee). We define social action as any action that
537 threatens or affirms the addressee's face. So, affirming and threatening face are social analogs to doing
538 moral good and harm respectively. In contexts where it is helpful, this definition also allows us to refer
539 to robots with different capacities to affect face as having different degrees of social agency, rather than
540 viewing social agency as a strictly binary attribute. We also propose that the term "social *actor*" can refer
541 to interactive entities capable of social action, but lacking the other criteria for agency (autonomy and/or
542 adaptability).

543 Some scholars have opined that it is common to view social agents as equivalent to "communicating
544 agents" (Castelfranchi, 1998), and thus might simply say that any communicative action is a social action.

545 Though the ability to nontrivially communicate implies the capacity to threaten face, we choose to base our
546 definition of social action directly on face because it allows for a more intuitive parallel to moral agency
547 without excluding any meaningful communicative actions. The vast majority of communicative actions
548 that an agent can perform have the capacity to impact face. Just in terms of face threat, any kind of request,
549 reminder, warning, advice, offer, commitment, compliment, or expression of negative emotion threatens
550 the addressee's negative face, and any criticism, rebuke, insult, disagreement, irreverence, boasting, non-
551 cooperation, or raising of divisive topics threatens the addressee's positive face (Brown and Levinson,
552 1987). A single speech act can carry several elements that affect face in different ways, and even the mere
553 act of purposefully addressing someone is slightly affirming of their positive face by acknowledging them
554 as worth addressing, and slightly threatening of their negative face by imposing on their time. Indeed, it is
555 difficult to think of a meaningful communicative action that would have no impact on face.

556 Another reason to ground social action in face is because face is more concrete and computationalizable
557 than some other options (e.g., induced perceptions of human likeness or influence on emotional state),
558 while still being broad enough to encompass the whole set of actions that we would intuitively consider
559 to be social. There exist various parameterizations or pseudo-quantifications of face threat/affirmation,
560 including Brown and Levinson's own formula which presents the weight of a face threatening act (W)
561 as the sum: $W = D(S, H) + P(H, S) + R$ where $D(S, H)$ is the social distance between the speaker
562 (S) and hearer (H), $P(H, S)$ quantifies the power that H has over S , and R represents the culturally and
563 situationally defined level of imposition that the face threatening act entails. For negative face threatening
564 acts, R includes the expenditure of time and resources. For positive face threatening acts, R is harder to
565 determine, but it is given by the discrepancy between H 's own desired self-image and that presented in the
566 face threatening act. Individual roles, obligations, preferences, and other idiosyncrasies are subsumed into
567 R . Of course, the constituent parts of this equation cannot be precisely quantified in any canonical way
568 (nor can, for example, influence on behavioral or emotional status). We do not view this as a weakness
569 because we would not expect to precisely quantify the magnitude of socialness in an action. Humans cannot
570 precisely answer questions like "How social is it to hug your grandmother?" or "Which is more social,
571 asking a stranger for the time or tipping your waitress?". However, this equation nonetheless illustrates
572 some of the concrete underpinnings of face and shows how face connects to concepts like relational power,
573 interpersonal relationships, material dependence, cultural mores, etc.

574 Robots are valid sources of social action under this face-based definition. Typical task-oriented paradigms
575 of HRI involve robots either accepting or rejecting human requests (which either affirms or threatens both
576 positive and negative face), or making requests of humans (which threatens negative face). Even simply
577 informing human teammates about the environment threatens negative face by implying that the humans
578 ought to act based on the new information. Less task-oriented cases, like companionship robots for the
579 elderly (Heerink et al., 2010), also require face affecting social actions, though these may tend to be more
580 face affirming than in task-based interaction. Again taking the SnackBot Lee et al. (2012) as an example,
581 bringing someone a requested snack is face affirming, and so are dialogue behaviors like complimenting
582 snack choice or apologizing for delays. The SnackBot's dialogue behavior of asking people to move out of
583 the way is face threatening. Research examining how robots influence human face and how humans react
584 to robotic face threatening actions is ongoing (Jackson et al., 2019, 2020).

585 In comparison to our definition, Castelfranchi (1998) define an action as either social or nonsocial
586 depending on its purposive effects and the mind of the actor. Their social actions must be *goal-oriented* and
587 motivated by *beliefs* about predicted effects in relation to some goal. Their social actions are mainly based
588 on some exercise of power, to attempt to influence the behavior of other agents by changing their minds.

589 They specifically say that social action cannot be a behavioral notion based solely on external description.
590 This definition is not well-suited to our purposes because these internal underpinnings are unknowable
591 to a typical robot user, and thus preclude the user from viewing a robot as a social agent. We saw similar
592 reasoning in our decision to exclude goal-orientedness as a prerequisite for agency. Even if a user chooses
593 to adopt an intentional stance (see (Dennett, 1978)) toward a robot and infer goals motivating its behavior,
594 this does not imply that the robot actually has an internal representation of a goal or of the intended effects
595 of its actions; the person's intentional stance would only allow them to take social action towards the robot,
596 not vice versa. Given the popular perception of robots as social and the academic tendency to call them
597 social agents, we do not want a definition of social action that cannot apply to robot action or that relies on
598 factors that cannot be observed from a user's LoA. Furthermore, Castelfranchi's definition excludes, for
599 example, end-to-end deep neural dialogue systems that may not explicitly represent goals, beliefs, causality,
600 or interactants as potential sources of social action, but whose actions can clearly come across as social and
601 carry all the corresponding externalities. Our face-based definition does not have these limitations.

602 To be clear, our decision to define social action via face is not an arbitrary design choice, but rather a
603 result of face's integral role in all social interaction. We believe that an action's relationship to face is,
604 unavoidably and fundamentally, what determines whether that action is social because face is what creates
605 the experience of having social needs/desires in humans. It follows that, for robots, the appearance or
606 attribution of face, or some relationship to others' face, is what allows them to be social actors. Any action
607 that affects face is necessarily social, and any action that does not is necessarily asocial. This aligns well
608 with widespread intuitions about sociality and common parlance use of the term.

609 **2.3 Social Patency as Having Face**

610 Any social action must have a recipient whose face is affected. If social agency is an agent's capacity to
611 be a source of social action (to affirm or threaten face), then the corresponding notion of social *patency* is
612 the capacity to have one's face threatened or affirmed (i.e., having face). This is similar to the notion of
613 moral patency as the capacity to be benefited or harmed by moral action. Clearly, conscious humans are
614 simultaneously moral and social agents and patients at any reasonable LoA. However, neither moral nor
615 social patency at any given LoA strictly requires moral or social agency at the same LoA, which leads us
616 to the question of whether our robotic moral/social agents in HRI are also moral/social patients.

617 It seems clear that, at a reasonable LoA for a human interactant, it is possible to harm a robot,
618 making the robot a moral patient. This is especially clear for robots capable of affective displays of
619 protest and distress (Briggs and Scheutz, 2014). Indeed people deliberately abuse robots with surprising
620 frequency (Nomura et al., 2015). However, at a deeper LoA, we know that current robots cannot feel pain
621 (or pleasure), have no true internal emotional response to harm like fear, and lack the will towards self
622 preservation inherent in most lifeforms. Thus, at this deeper LoA the robot is not a moral patient.

623 Likewise, a robot's social patency depends on the LoA considered. It is feasible to program a robot to
624 manifest behaviors indicating face wants, like responding negatively to insults and positively to praise, in
625 which case it would be a social patient at the user's LoA. However, at the developer's LoA, the robot still
626 has no face.

627 **2.4 Social and Moral Agencies as Independent**

628 We now discuss the extent to which social agency and moral agency can manifest in machines independent
629 of one another. We believe that some machines, including some robots, are largely perceived as asocial
630 moral agents, while others are seen as amoral social agents. Although, for the most part, social robots
631 do not fall in either of these groups, we believe that they are worth presenting as points of reference for

632 understanding the special moral and social niche occupied by language capable robots. We continue to
633 consider these technologies from the user's LoA.

634 Some artificial agents are popularly ascribed some form of moral agency without behaving socially or
635 even possessing the capacity for communication outside of a narrow task-based scope. We call such agents
636 "asocial moral agents", and use autonomous motor vehicles as the quintessential example. If we include the
637 likely possibility that autonomous vehicles will learn and change their behavior in response to changing
638 road conditions or passenger preferences, they are agentic at the passenger's LoA by being interactive,
639 autonomous, and adaptive.

640 In terms of moral action, while autonomous motor vehicles are obligated to conform to the legal rules
641 of the road, they are also expected to engage in extralegal moral decision making and moral reasoning.
642 Myriad articles, both in popular culture and in academia, contemplate whether and how autonomous cars
643 should make decisions based on moral principles (e.g., (Bonneton et al., 2016)). Questions like "in an
644 accident, should the car hit a school bus to save its own passenger's life? Or should it hit the barrier and kill
645 its passenger to save the school children?" have taken hold of popular imagination and proliferated wildly.
646 Regardless of the actual usefulness of such questions (cf. Himmelreich (2018)), it is clear that autonomous
647 cars are being ascribed moral agency.

648 We can also consider whether autonomous vehicles might be capable of social action. For example, using
649 a turn signal is clearly communicative, but it is also legally mandated; an autonomous vehicle would signal
650 an impending turn regardless of whether any other driver was present to see the turn signal. Given the legal
651 motivation behind the turn signal and the fact that it has no specific intended addressee, we view it as the
652 rare communicative act with no (or negligible) impact to face. Indeed, any communication via turn signal
653 would be considered incidental to law-following by the typical driver. Other driving behavior can also be
654 communicative; though we do not expect autonomous vehicles to engage in tailgating or road rage, we
655 could imagine that they might change the norms governing human driving behavior by modeling those
656 norms themselves. For example, if all autonomous vehicles on the road adopt a uniform following distance,
657 this behavior might influence human drivers sharing the road to do the same. However, this potential
658 normative influence is distinct from that of social robots in that it is passive, incidental, unintentional, and
659 not principally communicative, and therefore not face-relevant.

660 In other cases, depending on behavior, robots could be perceived as amoral social agents. Social robots
661 that do not have the ability to act on their environment in any meaningful extra-communicative capacity may
662 be physically unable (or barely able) to produce moral action. As an example, consider MIT's Kismet robot,
663 which is expressive, (non-linguistically) communicative, and social, but largely helpless and incapable
664 of acting extra-communicatively. Many social actions are available to Kismet. For example, making a
665 happy expression/noise when a person enters the room is face affirming, and a disgusted expression face
666 threatening. Given the right behaviors, Kismet could also meet our prerequisites for agency and be an
667 amoral social agent.

668 When moral and social agency are both present, as is the case for most social robots at the user's LoA,
669 their combination gives rise to interesting phenomena. Social robots can occupy a unique sociotechnical
670 niche: part technological tool, part agentic community member. This status allows robots to play an active
671 role in shaping the community norms that inform human morality, which behavioral ethics has shown
672 to be dynamic and malleable (Gino, 2015). And while robots are not the only technology to play a role
673 in shaping human norms (Verbeek, 2011), we believe their social agency grants them uniquely powerful
674 normative influence. For example, robots have been shown to hold measurable persuasive capacity over

675 humans, both via explicit and implicit persuasion (Briggs and Scheutz, 2014; Kennedy et al., 2014), and
676 even to weaken human (application of) moral norms via simple question asking behavior (Jackson and
677 Williams, 2019).

678 Language capable robots are unique among technologies not only in the strength of their potential moral
679 influence, but also in their ability to take an active and purposeful role in shaping human moral norms (or
680 human application of moral norms) as social agents. However, this capability is a double-edged sword. On
681 the one hand, robots of the future could productively influence the human moral ecosystem by reinforcing
682 desirable norms and dissuading norm violations. On the other hand, today's imperfect moral reasoning
683 and natural language dialogue systems open the door for robots to inadvertently and detrimentally impact
684 the human moral ecosystem through reasoning errors, miscommunications, and unintended implicatures.
685 It is thus crucial to ensure moral communication and proper communication of moral reasoning from
686 robots, especially in morally consequential contexts. The power to transfer or alter norms comes with the
687 responsibility to do so in a morally sensitive manner.

3 REVISITING RELATED WORK

688 Revisiting the theories of social agency from Section 1.2, we see that our definition is more inclusive than
689 that of Nagao and Takeuchi (1994) and Alač (2016) in that we demphasize the robot's embodiment and
690 materiality to account for purely digital potential social agents that we see in HRI research (Lee et al., 2006;
691 Heerink et al., 2010), and do away with the teleological and internal considerations (e.g., goal-orientedness
692 and intentionality) that would not be knowable to the typical robot user (cp. (Levin et al., 2013; Pollini,
693 2009)). On the other hand, our work is more restrictive than Pollini (2009) because we exclude "entities by
694 imagination" as potential social agents, and specify that there are several behavioral traits necessary for
695 social agency. This approach balances the more human-ascription-centered and more robot-trait-centered
696 conceptualizations of social agency. Our theory acknowledges the human role in determining social agency
697 by centering human face and the human's LoA, without reducing social agency to the mere ascription
698 thereof. At the same time, we concretely describe the robot traits necessary for social agency at a given
699 LoA.

700 Revisiting the studies from Section 1.3, which referenced social agents and social agency without
701 principally focusing on defining those concepts, we see that our definition can encompass the wide diversity
702 of potential social agents in HRI. Particularly at the user's LoA, robots can be social agents regardless of
703 embodiment, teleoperation, task-orientedness, morphology, mobility, or linguistic capacity. However, some
704 of the robots we reviewed would actually be excluded by our definition at the user's LoA by failing to
705 meet behavioral prerequisites, particularly by lacking indications of adaptability (e.g., (Lee et al., 2006;
706 Roubroeks et al., 2011; Heerink et al., 2010)). Interestingly, robots with a human teleoperator, like the
707 SnackBot (Lee et al., 2012) might be *more* likely to be socially agentic at the user's LoA than those with
708 simpler self-controlled behavior.

709 Finally, we stress that our theory complements (rather than competes with) much of the previous work
710 we discussed. For example, some of the proxemic and haptic human behavior that Alač (2016) observed in
711 their ethnographic study, like the choice to touch a robot's forearm rather than other body parts, might be
712 understood within our theory as stemming from attributions of social *patience* to the robot, rather than
713 social agency. Likewise, our conception of social agency may well be tied to, for example, psychological
714 reactance (Roubroeks et al., 2011) or trust (Ullman et al., 2014).

4 CONCLUDING REMARKS

715 We have presented a theory of social agency wherein a social agent (a thing with social agency) is any
716 *agent* capable of *social action* at the *LoA* being considered. A *LoA* is a set of observables, and the *LoAs*
717 most relevant to our discussion have been the robot user's, the developer's (or system *LoA*), and, to a lesser
718 extent, the architecture *LoA*. *Agency* at any given *LoA* is determined by three criteria which we defined
719 concretely above: interactivity, autonomy, and adaptability. We have defined *social action* as any action
720 that threatens or affirms the addressee's face, and refer to the addressee in this scenario as a social patient.
721 More specifically, *social patiency* is the capacity to be the recipient of social action, i.e., having face. These
722 definitions came from parallel concepts in the philosophy of *moral agency* (Floridi and Sanders, 2004).
723 We motivated our theory of social agency by presenting a sample of the inconsistent, underspecified, and
724 problematic theories and usages of social agency in the HRI literature.

725 Based on our theory, we have several recommendations for the HRI community. We recognize a tendency
726 to casually use the word "agent" to refer to anything with any behavior, and to correspondingly use "social
727 agent" to simply mean "social thing." We encourage authors to consider either switching to the broader term
728 "social actor" as defined above, or to briefly specify that they are using the term "social agent" informally
729 and do not intend to imply social agency in any rigorous sense. We further recommend that any paper
730 dealing with social agency be specific in selecting a suitable definition (such as the one presented in this
731 work) and *LoA*.

732 It will be important for future studies to develop, refine, and validate measurements of social (and
733 moral) agency. There exists early work on developing a survey to measure "perceived moral agency" for
734 HRI (Banks, 2019), however some questions seem to conflate moral *goodness* with moral *agency*, and,
735 despite measuring facets of autonomy and moral *cognition*, the survey does not measure the capacity
736 for taking moral *action*. Some of the proxies that we saw used for social agency in Section 1.3, like
737 human-likeness, realness, and livingness (Ghazali et al., 2019) do not match our new conceptualization of
738 social agency. Others, like gaze (Baxter et al., 2014), could be promising but have yet to be validated with
739 our theory (or, to our knowledge, any particular theory) of social agency in mind. Validated metrics would
740 facilitate experimental work motivated by our theory.

741 For example, future work designed to evaluate and further concretize our theory could empirically verify
742 whether changing the *LoA* at which somebody is viewing a robot causes a corresponding change to their
743 assessment of that robot as a (social) agent. The results could either strengthen the argument that the *LoA*
744 is a critical prerequisite for the discussion of agency, or indicate that colloquial conceptions of agency do
745 not account for *LoA*, despite its importance in rigorous academic discussions. Another avenue for this type
746 of work would be to manipulate the magnitude of face threat/affirmation that a social robot is capable of
747 and examine how that manipulation effects perceptions of the robot as a social agent. This experiment
748 would specifically target our definition of social action as grounded in face.

749 Measures of social agency would also allow us to examine its relationship with persuasion and trust. On
750 the one hand, we could imagine that decreasing a robot's social agency (by lowering its propensity to affect
751 face) could increase its persuasive capacity if people are more amenable to persuasion when their face is
752 not threatened. On the other hand, increasing a robot's social agency might increase its persuasive capacity
753 if people are more likely to trust a more human-like robot.

754 Furthermore, it will be important to probe for causal relationships between ascriptions of social agency
755 and ascriptions of moral responsibility and competence in robots. In human children, development of
756 increased capacity for social action is typically correlated with development of other facets of intelligence

757 and skills, including moral reasoning. However, this correlation does not necessarily exist for robots, since
758 a robot could be socially agentic and competent, with a wide range of possible social actions, and still have
759 no moral reasoning capacity. If robot social agency, or social behavior in general, leads interactants to
760 assumptions of moral competence or overall intelligence (as it likely would in humans), this could lead
761 to dangerous overtrust in robot teammates in morally consequential contexts that they are not equipped
762 to handle. Thus, giving a robot linguistic/social competence would also necessitate giving the robot a
763 corresponding degree of moral competence.

764 Finally, though there is evidence for an ontological distinction between humans and robots (Kahn et al.,
765 2011), it is not yet clear where differences (and similarities) will manifest in terms of moral and social
766 agency. We will require human points of reference in future HRI studies to fully understand how the
767 emerging moral and social agency of robots relate to those qualities in humans.

CONFLICT OF INTEREST STATEMENT

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

AUTHOR CONTRIBUTIONS

770 This paper was written by RBJ with input, feedback, and review by TW.

FUNDING

771 This work was funded in part by Air Force Young Investigator Award 19RT0497.

ACKNOWLEDGMENTS

772 We would like to thank Dylan Doyle-Burke for his collaboration during earlier stages of this work.

REFERENCES

- 773 Alač, M. (2016). Social robots: Things or agents? *AI & society* 31, 519–535
- 774 Atkinson, R. K., Mayer, R. E., and Merrill, M. M. (2005). Fostering social agency in multimedia learning:
775 Examining the impact of an animated agent's voice. *Contemporary Educational Psychology* 30, 117–139
- 776 Banks, J. (2019). A perceived moral agency scale: Development and validation of a metric for humans and
777 social machines. *Computers in Human Behavior* 90, 363–371
- 778 Baxter, P., Kennedy, J., Vollmer, A.-L., de Greeff, J., and Belpaeme, T. (2014). Tracking gaze over time in
779 hri as a proxy for engagement and attribution of social agency. In *Proceedings of the 2014 ACM/IEEE*
780 *international conference on Human-robot interaction*. 126–127
- 781 Billett, S. (2008). Learning throughout working life: a relational interdependence between personal and
782 social agency. *British Journal of educational studies* 56, 39–58
- 783 Bonnefon, J.-F., Shariff, A., and Rahwan, I. (2016). The social dilemma of autonomous vehicles. *Science*
784 352, 1573–1576
- 785 Briggs, G. and Scheutz, M. (2014). How robots can affect human behavior: Investigating the effects of
786 robotic displays of protest and distress. *Int'l Journal of Social Robotics*
- 787 Brown, P. and Levinson, S. (1987). *Politeness: Some Universals in Language Usage* (Cambridge University
788 Press)
- 789 Castelfranchi, C. (1998). Modelling social action for ai agents. *Artificial intelligence* 103, 157–182
- 790 Danaher, J. (2020). Welcoming robots into the moral circle: a defence of ethical behaviourism. *Science*
791 *and Engineering Ethics* 26, 2023–2049

- 792 Dennett, D. C. (1978). Three kinds of intentional psychology. *Perspectives in the philosophy of language: A concise anthology*, 163–186
- 793
- 794 Dobres, M.-A. and Hoffman, C. R. (1994). Social agency and the dynamics of prehistoric technology. *Journal of archaeological method and theory* 1, 211–258
- 795
- 796 Floridi, L. (2008). The method of levels of abstraction. *Minds and machines* 18, 303–329
- 797 Floridi, L. and Sanders, J. W. (2004). On the morality of artificial agents. *Minds and machines* 14, 349–379
- 798 Gardner, A. (2016). *Agency uncovered: Archaeological perspectives on social agency, power, and being human* (Routledge)
- 799
- 800 Garibay, J. C. (2015). Stem students' social agency and views on working for social change: Are stem disciplines developing socially and civically responsible students? *Journal of Research in Science Teaching* 52, 610–632
- 801
- 802
- 803 Garibay, J. C. (2018). Beyond traditional measures of stem success: Long-term predictors of social agency and conducting research for social change. *Research in Higher Education* 59, 349–381
- 804
- 805 Ghazali, A. S., Ham, J., Barakova, E., and Markopoulos, P. (2018). The influence of social cues in persuasive social robots on psychological reactance and compliance. *Computers in Human Behavior* 87, 58–65
- 806
- 807
- 808 Ghazali, A. S., Ham, J., Markopoulos, P., and Barakova, E. I. (2019). Investigating the effect of social cues on social agency judgement. In *HRI*. 586–587
- 809
- 810 Gino, F. (2015). Understanding ordinary unethical behavior: Why people who value morality act immorally. *Current opinion in behavioral sciences* 3, 107–111
- 811
- 812 Gray, K. and Wegner, D. M. (2009). Moral typecasting: divergent perceptions of moral agents and moral patients. *Journal of personality and social psychology* 96, 505
- 813
- 814 Heerink, M., Kröse, B., Evers, V., and Wielinga, B. (2010). Assessing acceptance of assistive social agent technology by older adults: the almere model. *International journal of social robotics* 2, 361–375
- 815
- 816 Himmelreich, J. (2018). Never mind the trolley: The ethics of autonomous vehicles in mundane situations. *Ethical Theory and Moral Practice* 21, 669–684
- 817
- 818 Jackson, R. B., Wen, R., and Williams, T. (2019). Tact in noncompliance: The need for pragmatically apt responses to unethical commands. In *AAAI Conference on Artificial Intelligence, Ethics, and Society*
- 819
- 820 Jackson, R. B. and Williams, T. (2019). Language-capable robots may inadvertently weaken human moral norms. In *Companion Proceedings of the ACM/IEEE International Conference on Human-Robot Interaction (alt.HRI)*
- 821
- 822
- 823 Jackson, R. B., Williams, T., and Smith, N. (2020). Exploring the role of gender in perceptions of robotic noncompliance. In *Proceedings of the 2020 ACM/IEEE International Conference on Human-Robot Interaction*. 559–567
- 824
- 825
- 826 Johnson, D. G. and Miller, K. W. (2008). Un-making artificial moral agents. *Ethics and Information Technology* 10, 123–133
- 827
- 828 Kahn, P. H., Reichert, A. L., Gary, H. E., Kanda, T., Ishiguro, H., Shen, S., et al. (2011). The new ontological category hypothesis in human-robot interaction. In *2011 6th ACM/IEEE International Conference on Human-Robot Interaction (HRI) (IEEE)*, 159–160
- 829
- 830
- 831 Kennedy, J., Baxter, P., and Belpaeme, T. (2014). Children comply with a robot's indirect requests. In *HRI*
- 832
- 833 Lee, K. M., Jung, Y., Kim, J., and Kim, S. R. (2006). 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, 962–973
- 834
- 835 Lee, M. K., Kiesler, S., Forlizzi, J., and Rybski, P. (2012). Ripple effects of an embedded social agent: a field study of a social robot in the workplace. In *Proceedings of the SIGCHI Conference on Human*
- 836

- 837 *Factors in Computing Systems*. 695–704
- 838 Levin, D. T., Adams, J. A., Saylor, M. M., and Biswas, G. (2013). A transition model for cognitions about
839 agency. In *2013 8th ACM/IEEE International Conference on Human-Robot Interaction (HRI)* (IEEE),
840 373–380
- 841 Luria, M., Hoffman, G., Megidish, B., Zuckerman, O., and Park, S. (2016). Designing vyo, a robotic smart
842 home assistant: Bridging the gap between device and social agent. In *2016 25th IEEE International
843 Symposium on Robot and Human Interactive Communication (RO-MAN)* (IEEE), 1019–1025
- 844 Marchesi, S., Ghiglino, D., Ciardo, F., Perez-Osorio, J., Baykara, E., and Wykowska, A. (2019). Do we
845 adopt the intentional stance toward humanoid robots? *Frontiers in psychology* 10, 450
- 846 Meyer, J. W. and Jepperson, R. L. (2000). The ‘actors’ of modern society: The cultural construction of
847 social agency. *Sociological theory* 18, 100–120
- 848 Nagao, K. and Takeuchi, A. (1994). Social interaction: Multimodal conversation with social agents. In
849 *AAAI*. vol. 94, 22–28
- 850 Nass, C., Steuer, J., and Tauber, E. R. (1994). Computers are social actors. In *Proceedings of the SIGCHI
851 conference on Human factors in computing systems* (ACM), 72–78
- 852 Nomura, T., Uratani, T., Kanda, T., Matsumoto, K., Kidokoro, H., Suehiro, Y., et al. (2015). Why do
853 children abuse robots? In *HRI Extended Abstracts*. 63–64
- 854 Perez-Osorio, J. and Wykowska, A. (2019). Adopting the intentional stance towards humanoid robots. In
855 *Wording Robotics* (Springer). 119–136
- 856 Pollini, A. (2009). A theoretical perspective on social agency. *AI & society* 24, 165–171
- 857 Roubroeks, M., Ham, J., and Midden, C. (2011). When artificial social agents try to persuade people:
858 The role of social agency on the occurrence of psychological reactance. *International Journal of Social
859 Robotics* 3, 155–165
- 860 Schellen, E. and Wykowska, A. (2019). Intentional mindset toward robots—open questions and
861 methodological challenges. *Frontiers in Robotics and AI* 5, 139
- 862 Straub, I. (2016). ‘it looks like a human!’ the interrelation of social presence, interaction and agency
863 ascription: a case study about the effects of an android robot on social agency ascription. *AI & society*
864 31, 553–571
- 865 Thellman, S., Silvervarg, A., and Ziemke, T. (2017). Folk-psychological interpretation of human vs.
866 humanoid robot behavior: Exploring the intentional stance toward robots. *Frontiers in psychology* 8,
867 1962
- 868 Thellman, S. and Ziemke, T. (2019). The intentional stance toward robots: Conceptual and methodological
869 considerations. In *The 41st Annual Conference of the Cognitive Science Society, July 24-26, Montreal,
870 Canada*. 1097–1103
- 871 Ullman, D., Leite, L., Phillips, J., Kim-Cohen, J., and Scassellati, B. (2014). Smart human, smarter
872 robot: How cheating affects perceptions of social agency. In *Proceedings of the Annual Meeting of the
873 Cognitive Science Society*. vol. 36
- 874 Verbeek, P.-P. (2011). *Moralizing Technology: Understanding and Designing the Morality of Things*
875 (University of Chicago Press)
- 876 Westlund, J. M. K., Martinez, M., Archie, M., Das, M., and Breazeal, C. (2016). Effects of framing a
877 robot as a social agent or as a machine on children’s social behavior. In *2016 25th IEEE International
878 Symposium on Robot and Human Interactive Communication (RO-MAN)* (IEEE), 688–693