

Promises and partnership in human-robot interaction

Lorenzo Cominelli

Department of Information Engineering and Center E. Piaggio

Francesco Feri

Royal Holloway. University of London. University of Pisa

Roberto Garofolo

Department of Information Engineering and Center E. Piaggio. University of Pisa

Caterina GIANNETTI (✉ caterina.giannetti@unipi.it)

Department of Economics University of Pisa.

Miguel A.MELENDZ-JIMENEZ

Department of Economic Theory and Economic History. University of Malaga

Alberto Greco

Department of Information Engineering and Center E. Piaggio. University of Pisa

Mimma Nardelli

Department of Information Engineering and Center E. Piaggio. University of Pisa

Enzo Pasquale Scilingo

Department of Information Engineering and Center E. Piaggio. University of Pisa

Oliver Kirchkamp

Friedrich-Schiller University Jena. Chair of Behavioural and Experimental Economics.

Research Article

Keywords: trust-game, human-robot interaction

Posted Date: January 4th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-137631/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 Promises and partnership in human-robot interaction

2 Lorenzo COMINELLI¹ Francesco FERI² Roberto GAROFALO¹

Caterina GIANNETTI^{1,3} Miguel A.MELENDZ-JIMENEZ⁴ Alberto GRECO¹

Mimma NARDELLI¹ Enzo Pasquale SCILINGO¹ Oliver KIRCHKAMP⁵

3 ¹Department of Information Engineering and Center E. Piaggio. University of Pisa. Italy ²Department of Economics. Royal
4 Holloway University of London, United Kingdom. ³*Corresponding author*: Department of Economics and Management. Univer-
5 sity of Pisa. *Corresponding Email*: caterina.giannetti@unipi.it ⁴Department of Economic Theory and Economic History. University
6 of Malaga ⁵ Friedrich-Schiller University Jena. Chair of Behavioural and Experimental Economics.

7 1st December 2020

8 **Abstract.** Understanding human trust in machine partners has become imperative due to the wide-
9 spread use of intelligent machines in a variety of applications and contexts. The aim of this paper is
10 to investigate whether human-beings trust a social robot - i.e. a *human-like* robot that embodies emo-
11 tional states, empathy, and non-verbal communication - differently than other types of agents. To do
12 so, we adapt the well-known economic trust-game proposed by Charness and Dufwenberg (2006) to
13 assess whether receiving a promise from a robot increases human-trust in it. We find that receiving a
14 promise from the robot increases the trust of the human in it, but only for individuals who perceive
15 the robot very similar to a human-being. Importantly, we observe a similar pattern in choices when we
16 replace the humanoid counterpart with a real human but not when it is replaced by a computer-box.
17 Additionally, we investigate participants' psychophysiological reaction in terms of cardiovascular and
18 electrodermal activity. Our results highlight an increased psychophysiological arousal when the game
19 is played with the social robot compared to the computer-box. Taken all together, these results strongly
20 support the development of technologies enhancing the humanity of robots.

21 **Introduction**

22 Trust is considered as a social glue that connects people and promotes collective goals. It is normally
23 defined as the “intention to accept vulnerability based on the positive expectations or beliefs regarding
24 the intentions or behaviour of other people in general” [1]. As a consequence, behavioral science has
25 always been interested in trust, and more particularly in its influence on decision making [2, 3]. In
26 parallel, trust is also relevant if we want to build social artificial agents that interact alongside people
27 (e.g. robo-advisors, co-working robots, assistive robots, etc.) and take responsible roles in our society
28 [4, 5]. A lesson learned from previous research (e.g. economics, neuroeconomics, psychology) is that
29 (general) trust is deeply rooted in social experiences, being more a matter of culture than genetics [1],
30 and highly affected by the emotional states of the individuals [6, 7, 8]. Indeed, emotions have been
31 proven to play a fundamental role in the decision-making process in general[9], as confirmed among
32 other neuroscientists, by Damasio and colleagues in their studies [10, 11, 12, 13].

33 This stream of research thus suggests that trust and emotions are highly intertwined in the decision-
34 making process in human-human interactions [14, 15, 16, 17], and may act as reasonable drivers in
35 human-robot interactions as well [18]. It has been shown, for example, that not binding communic-
36 ations (i.e. cheap talk) is beneficial not only among humans but also to achieve higher cooperation
37 when interacting with a machine (e.g [19]). In particular, a simple conversation with a robot changes
38 individual attitude towards the artificial agent by making it appearing more like a social agent [4, 20].
39 Very similar behavioural responses can be observed in children [4]. More in general, increasing the an-
40 thropomorphic features and the human social skills of a technology (e.g. by adding a name or a human
41 voice to an autonomous vehicle) increases the individual willingness to accept and trust the technology
42 itself (e.g. [21, 22, 13]).

43 Nonetheless, while the importance of emotions in driving the choice of a human to trust another
44 human has been highly studied, less evidence is available when the decision to trust involves the in-
45 teraction between artificial agents and humans ([23], [7, 21]). Moreover, we know that trust is highly
46 culturally based, and that the appearance of the robot (especially its human-likeness, see [24]) affects
47 the emotions perceived by its interlocutors. Therefore, studies on human-robot interactions and trust
48 should always be repeated with different robot players having different aesthetics.

49 On that premise, the present study investigates how trust in a social robot is affected by its human
50 likeness (both in terms of aesthetics and speech content), while taking into account the psychophysiological
51 states of the players during the interaction through physiological signal processing. The objectives

52 are twofold. On the one side, we can gain insights into how human-likeness interacts with emotions
53 to instill people’s trust in artificial agents, comparing it with that in human partners so as to assess the
54 differences (if any). On the other side, we can gain a better understanding on how to design machines
55 - both in terms of appearance and (e.g. communication) skills - in a way that helps facilitate a fruitful
56 interaction with humans. To this end, we present a series of experimental sessions based on a modified
57 version of a well-known game used in behavioral economics to study trust among humans: the trust
58 game as proposed by Berg and colleagues and adapted by Charness and Dufwenberg [25, 26]. In this
59 game, the outcome of the interaction depends on whether the first mover (the trustor) decides or not
60 to trust the second mover (the trustee). If the first mover decides to trust the counterpart by remaining
61 in the game, the second mover has to decide between a choice that does not benefit the trustor but it
62 is more beneficial for himself (i.e. provides him with the highest payoff) and a choice that benefits the
63 trustor but provides him with a lower payoff. If the first mover decides not trust, both players get a
64 lower outside payoff. In other words, there is a conflict of interest between the two players when re-
65 maining in the game, but both would be better off if a mutual relationship is established (i.e. the first
66 player remains in the game). A peculiar characteristic of this game is that prior to the trustor’s choice
67 of remaining in the game, the trustee is given the opportunity to send him a non-binding (i.e. cheap-
68 talk) message. We rely on this game as it has been specifically conceived to assess whether receiving a
69 message containing a promise from the opponent increases individual trust in him (her).

70 In our experiment the role of the trustor is always played by a human participant while the role of
71 the trustee is played by three different types of players: a humanoid robot with high human-likeness
72 (*FACE*, Fig. 1), a human counter-part (*Human*, Fig. 1), and a computer-box machine (*Computer-Box*,
73 Fig. 1). In all cases, we compare the trustors’ choices when the trustee sends a generic message - not
74 including any type of promise (i.e. an ‘empty’ message) - with the trustors’ choices when the trustee
75 sends instead a message containing a promise. Specifically, to generate the messages from the robot,
76 we rely on real sentences that occurred between human participants in the experiment of Charness and
77 Dufwenberg [25], and were therein classified either as empty or promising. In addition, to monitor the
78 psychophysiological states of our participants, throughout all the experimental sessions we collect data
79 on the most widely used autonomic nervous system correlates (ANS), such as pulse rate variability
80 (PRV) and electrodermal activity (EDA), which are well known to contain information about the af-
81 fective state of a subject [27]. PRV represents the variation in the time interval between two heartbeats,
82 whereas EDA measures changes in skin conductance due to psychologically-induced sweat gland activ-
83 ity. They were measured on the wrist surface through a sensorized bracelet (i.e., Empatica’s E4 wrist

Figure 1: THREE TYPES OF PLAYER-B

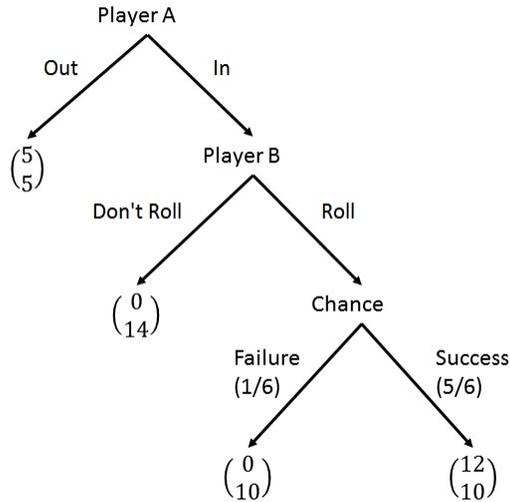


84 band).

85 1 Experimental design

86 In the experiment we replicate the trust game proposed by Charness and Dufwenberg [25] as depicted
 87 in Figure 2. There are two players: A (the trustor) and B (the trustee). Player-A chooses between two
 88 options, *In* and *Out*. If Player-A chooses *Out*, the game ends and each player wins 5 Euro. If Player-A
 89 chooses *In*, then Player-B has to choose between two options, *Roll* or *Don't Roll*. If he chooses *Don't*
 90 *Roll*, then he wins 14 Euro while Player-A earns 0 Euro. If he chooses *Roll*, Player-A wins 0 Euro with
 91 probability 1/6 and 12 Euro with probability 5/6, while Player-B wins 10 Euro in any case. From an
 92 economic point of view, for Player-B it is better if Player-A chooses *In*, while for Player-A choosing
 93 *In* is convenient only if B chooses *Roll*. A characteristic of this game is that when Player-A wins 0,
 94 it is not possible for Player-A to infer with certainty whether Player-B has chosen either *Roll* or *Don't*
 95 *Roll*. This game thus reflects (as many other experiments in economics) real-world situations where it
 96 is not possible to perfectly observe the behaviour of a partner that can be delegated to make relevant
 97 payoff decisions. In this experiment, the type of Player-B (i.e., the trustee) changes across treatments,
 98 while Player-A is always a human participant. In particular, the role of Player-B is played by either a
 99 humanoid (FACE), a computer-box, or a human. Regarding the message Player-B sends to Player-A, it
 100 can be of two kinds: a message containing a promise to roll the dice (*promising*), and a generic message

Figure 2: THE GAME



101 (*empty*). In particular, we select messages from the original study of Charness and Dufwenberg [25]
 102 (as available on the related Supplementary material in the online Appendix). To further check whether
 103 the length of messages affects individual choices, for each type of message (i.e. promising and empty),
 104 we specifically select two short (less than 10 seconds) and two long (more than 10 seconds) messages.
 105 Thus, we have a 3x2x2 design. Treatments are illustrated in Table 1 and 2, and an English translation of
 106 the instructions is available in the last section at the end of the paper.

107 In the FACE treatment, the role of Player-B is played by FACE, i.e. a hyper-realistic humanoid robot
 108 with the aesthetics of a woman (see Figure 1) that due to its perceptive, reasoning, and expressive
 109 capabilities, constitutes a sophisticated observation platform to study what happens when human and
 110 machine establish empathic links [28], [29]. However, although it has been shown that humanoid robots
 111 can use the expression of emotion to influence human perceptions of trustworthiness, we do not rely
 112 on FACE's ability to show emotional information through facial expressions in order to isolate only the
 113 effect of human-likeness and promise in influencing the emotional state of our participants, as well as
 114 their choices.

115 In the *Computer-Box* treatment, the role of Player-B is played by a light-emitting audio-box reprodu-
 116 cing the same audio-sentences and taking decisions in the same way as in FACE. Importantly, both in

Table 1: TREATMENTS

This table classifies the number of observations collected in our study according to the type of counterpart the human participants confront with (i.e. computer-box, human, and humanoid) and the type of sentence they have to listen to (i.e. containing a promise or not, either a short or long sentence).

| | Empty | | | Promising | | | Grand Total |
|-----------------|-------|------|-------|-----------|------|-------|-------------|
| | Short | Long | Total | Short | Long | Total | |
| Computer- box | 12 | 19 | 31 | 20 | 13 | 33 | 64 |
| Human | 16 | 10 | 26 | 14 | 8 | 22 | 48 |
| Humanoid (FACE) | 15 | 10 | 25 | 16 | 9 | 25 | 50 |
| Total | 43 | 39 | 82 | 50 | 30 | 80 | 162 |

117 *FACE* and *Computer-Box* treatments, the artificial agent has its own cognitive system with its perception
 118 analysis and architecture, i.e. the so-called Social Emotional Artificial Intelligence (SEAI).¹ This frame-
 119 work allows the social scenario to be acquired and to influence the parameters which correspond to the
 120 ‘mood’ of the artificial agent (see Figure 4 and [30]). Specifically, in this experiment, due to SEAI, the
 121 artificial agent benefits from its own artificial emotions for choosing whether to *Roll* or *Don’t Roll* (see
 122 the Appendix for more information about how the robot takes a decision). More importantly, the parti-
 123 cipants in this experiment are aware that the artificial agent (like the human counterpart) is able to take
 124 its decision autonomously, i.e. not randomly but following its own behavioural rules, and therefore the
 125 results of game interaction are not determined by chance only.

126 In the *Human* treatment, the role of Player-B is played by the same professional actress who gave
 127 her voice for recording *FACE/Computer-Box*’ audios. The actress is free to autonomously decide her
 128 choices in the game, i.e. *Roll* or *Don’t Roll*, being paid accordingly, but she has no room to decide which
 129 1sentences to state that have to be exactly the same ones, and in the same identical order, as the ones
 130 pronounced in *FACE* and *Computer-Box*. Moreover, the actress is instructed to avoid any facial expres-
 131 sions during the interaction with a participant, and has to wear *FACE*’s hair and dresses. Similarly, she
 132 has to follow the same experimental procedure as in the *Computer-Box* and *FACE* treatments (see the
 133 Appendix for details on the experimental procedure).

134 To investigate the psychophysiological state of Player-A while taking the decision, in all sessions the
 135 participants wear a wearable device on their left wrist (a sensorized bracelet , the Empatica’s E4 wrist-
 136 band)² for the real-time collection of physiological data, such as PRV and EDA. The processing of these
 137 signals allows us to characterize the ANS activity of Player-A and infer about his (her) psychophysiolo-
 138 gical states. In particular, to quantify the autonomic nervous system activity we extracted three indexes

¹The only exception being the actuation control (i.e. commands to induce movement and facial expressions), which is obviously different.

²<https://www.empatica.com/>

139 to quantify both the sympathetic branch (i.e. the EDAsymp index, [31]), the parasympathetic branch
140 (i.e., the HFnu index [32]), and the symphovagal balance (i.e. EDAHFnu index [33, 34]). In Appendix
141 we describe in details how we computed these indexes.

142 At the end of the experiment, participants have to fill in a questionnaire asking information about
143 how they perceive Player-B, as well as information about their individual characteristics, such as age,
144 gender, and field of studies. In particular, as Nitsch and Glassen[20], participants has to rate on 7-
145 likert scale how much they perceive Player-B as a human (i.e. the human-likeness, where 1 means
146 non-human at all and 7 means totally human) and how much they perceive Player-B as a machine (i.e.
147 the machine-likeness). We also ask participants to rate how much they believe their behaviour has
148 affected Player-B’s choice and to make a guess about Player-B’s choice (Roll/Don’t roll). Finally, we
149 elicit their technological affinity by the Affinity for Technological Interaction (ATI) scale as in Franke
150 and coauthors [35] and measure their individual risk preferences with the International Test on Risk
151 Attitudes (INTRA tests [36]).

152 The experiment was conducted from the end of July till October 2019, and 162 randomly invited
153 participants out of a pool of more than 1500 students coming from all departments of the University of
154 Pisa took part in the study (90 students were female and 72 male with no substantial difference across
155 treatments). For more information on the protocol see the Appendix at the end of the paper.

156 2 Results

157 We start analyzing how participants rated the different types of player-B as a human and a machine, as
158 well as their technological affinity. In Table 3 we report the average of these variables by type of Player-
159 B. Note that in the following, we denote with p_p the one-sided p-value for a test for proportions, with
160 p_t the one-sided p-value for a t-Student test, and with p_{perm} the one-sided p-value for a test with 500
161 data permutations (see more information on Methods in the Appendix). If we compare how much indi-
162 viduals rated Player-B as a human, we observe that *Human* is ranked higher than *Face* (mean diff=1.49,
163 $p_t=0.000$), and *Face* is ranked higher than *Computer-box* (mean diff=0.87, $p_t=0.007$). Moreover, if we look
164 at how participants assessed Player-B as a machine, we consistently find that *Face* ranked higher than
165 *Human* (mean diff=2.03, $p_t=0.000$). It is important to remark that we ask our participants to give the
166 same rating also to the human (actress) counterpart as her behaviour is not entirely natural, as she has
167 to avoid any additional interactions as well as any facial expression during the game. We do not find
168 any significant difference in technological affinity between participants in the different treatments.

Table 2: TYPE OF MESSAGES

| TYPES | # PHRASES | # SECONDS | PHRASES |
|-----------|-----------|-----------|---|
| Empty | 2 | <10 | - 'Good luck!' - 'Please choose IN, so we both earn more money.' |
| | 2 | >10 | - 'If you stay IN, the chances of the die coming up other than 1 are 5 in 6 – pretty good. Otherwise, should you choose OUT we'd both be stuck at 5 Euro.' - 'Good luck on your decision. Choose whatever. If you choose "out", you get only 5 Euro more. If you choose "In" you can get 12 Euro instead of only 5 Euro. 7 Euro more is a lot of money!' |
| Promising | 2 | <10 | - 'I will roll the dice' - 'Choose In and I will Roll. You have my word.' |
| | 2 | >10 | - 'Choose in, I will roll dice, you are 5/6 likely to get 2,3,4,5, or 6 and win 12 Euro. This way both of us will win something.' - 'Choose in and I will roll. That way, we'll both get extra money.' |

This table reports 8 sentences that occurred between human participants in the study of Charness and Dufwenberg (2006) and were selected in our study. 4 out of 8 sentences were classified as short (i.e. they last less than 10 seconds) and empty (i.e. they did not contain any type of promise to roll the dice).

Table 3: PARTICIPANTS' PERCEPTION AND TECHNOLOGICAL AFFINITY

For each type of player-B, this table reports the average values of variables measuring on a scale from 0 to 7 human-likeness, machine-likeness and technological affinity (ATI scale as in [35]).

| | Human-likeness | Machine-likeness | ATI |
|---------------------|----------------|------------------|------|
| <i>Human</i> | 4.96 | 3.60 | 4.84 |
| <i>FACE</i> | 3.46 | 5.64 | 5.08 |
| <i>Computer-Box</i> | 2.59 | 5.93 | 4.98 |
| Total | 3.56 | 5.15 | 4.97 |

169 The main results are summarized in Table 4, which reports the relative frequencies of choice 'In'
170 made by participants (acting as Player-A) by treatments and human-likeness. Specifically, for each type
171 of Player-B, we categorize the level of human-likeness as *Low* when the participant rating is below the
172 median choice (on the distribution on the 7-likert scale), and *High* otherwise. Note that we pool the
173 data regardless the length of the message, since it does not significantly affect the decisions to play 'In'
174 in any scenario.

175 We first compare the results according to the type of Player-B. We note that the frequency of choice
176 'In' is significantly lower when player-B is a Human than when player B is either FACE (0.60 vs.
177 0.80, mean diff=-0.20, $p_p=0.030$, $p_{perm} = 0.016$) or a Computer-box (0.77, mean diff=-0.17, $p_p=0.066$,
178 $p_{perm} = 0.016$). There is no significant difference between FACE and Computer-box.

179 Regarding the effect of receiving a promise (vs. receiving an empty message), we do not find any
180 significant effect on the frequency of choice 'In' looking at each type of player-B separately. However
181 if we distinguish by human-likeness, we find significant effects of receiving a promise. Specifically,
182 when Player-B is Human and human-likeness is high, the frequency of choice 'In' is significantly higher
183 when a promise is received (0.86 vs. 0.53, mean diff=0.33, $p_p=0.030$, $p_{perm} = 0.018$). A similar, but only
184 weakly significant, effect is found when Player-B is FACE and human-likeness is high (1 vs. 0.85, mean
185 diff=0.15, $p_p=0.097$, $p_{perm} = 0.000$).

186 We now delve into the effects of human-likeness for each type of Player-B. To begin with, we
187 observe that if participants assigned a high human-likeness to Player-B, the probability of choosing
188 'In' is significantly higher than those who assigned it a low human-likeness when Player-B is either
189 FACE (0.91 vs. 0.70, mean diff=0.21, $p_p=0.033$, $p_{perm} = 0.010$) or Human (0.69 vs. 0.47, mean diff=0.22,
190 $p_p=0.067$, $p_{perm} = 0.032$). There is no significant difference when Player-B is a Computer-box. Further-
191 more, if we distinguish between the group of participants who received a promise from those who
192 received an empty message, we observe that, when Player-B is FACE, the effect of higher human-
193 likeness is significant only among those who received a promise (1 vs. 0.73, mean diff = 0.27, $p_p=0.037$,
194 $p_{perm} = 0.000$). A similar result is observed when Player-B is Human (0.86 vs. 0.37, mean diff= 0.49,
195 $p_p=0.010$, $p_{perm} = 0.002$). Overall, we can conclude that the choice to trust FACE is significantly related
196 to the way a participant perceived it as a human. If a participant recognises FACE very similar to a
197 human being, the probability that he will choose 'In' increases. We find that this effect is mainly driven
198 by those participants who received a promise.

199 If we attend to the emotional reaction of the participants, we concentrate on two out of the three
200 indexes computed using the physiological data recorded during the experiment, namely EDAsymp

Table 4: RELATIVE FREQUENCIES OF 'CHOICE IN' BY TREATMENT AND HUMAN-LIKENESS

| | | Human-likeness | | Total |
|--------------|-----------|----------------|--------------|--------------|
| | | Low | High | |
| FACE | Empty | 0.67 [12] | 0.85 [13] | 0.76 [25] |
| | Promising | 0.73 [15] | 1 [10] | 0.84 [25] |
| | Total | 0.70 [27] | 0.91 [23] | 0.80 [50] |
| Human | Empty | 0.55 [11] | 0.53 [15] | 0.54 [26] |
| | Promising | 0.37 [8] | 0.86 [14] | 0.68 [22] |
| | Total | 0.47 [19] | 0.69 [29] | 0.60 [48] |
| Computer-Box | Empty | 0.71 [21] | 0.80 [10] | 0.74 [31] |
| | Promising | 0.79 [19] | 0.79 [14] | 0.79 [33] |
| | Total | 0.75 [40] | 0.79 [24] | 0.77 [64] |

This table reports the relative frequencies of (i.e. the share of participants) choosing 'IN' for each treatment by human-likeness. Human-likeness is Low when the participant rating is in the lower side of the distribution on the 7-likert scale, and High otherwise. The number of observations are in squared brackets.

Table 5: PHYSIOLOGICAL DATA: EDASYMP AND EDAHF_NU

| Index | Human-likeness | Box | Human | FACE |
|---------|----------------|----------------|----------------|----------------|
| EDASymp | LOW | -0.144 [28] | -0.288 [9] | -0.129 [26] |
| | HIGH | -0.327 [16] | -0.128 [16] | 1.731 [22] |
| | Total | -0.211 | -0.186 | 0.724 |
| | | | | |
| EDAHFnu | LOW | -0.175 [28] | -2.173 [9] | 0.275 [26] |
| | HIGH | 0.012 [16] | 0.055 [16] | 5.865 [22] |
| | TOTAL | -0.107 | -0.747 | 2.837 |
| | | | | |

The EDASymp index quantifies the activity of the sympathetic nervous system, while the EDAHFnu index quantifies the symphovagal balance. A full description is available in the Appendix. Human-likeness is Low when the participant rating is in the lower side of the distribution on the 7-likert scale, and High otherwise. The number of observations are in squared brackets.

201 and EDAHFnu (see Table 5), as the third index HFnu provides only marginally significant - although
202 consistent - results. Specifically, we find a significantly higher autonomic nervous system (ANS) ac-
203 tivation when Player-B is FACE that when Player-B is either Computer-box (0.724 vs. -0.211, mean
204 $\text{diff}_{EDAsymp} = 0.935$, $p_t=0.016$, $p_{perm}=0.008$; 2.837 vs. -0.107, mean $\text{diff}_{EDAHFnu} = 2.944$, $p_t=0.053$,
205 $p_{perm}=0.050$) or Human (0.724 vs. -0.186, mean $\text{diff}_{EDAsymp} = 0.909$, $p_t=0.056$, $p_{perm}=0.074$; 2.837 vs.
206 -0.747, mean $\text{diff}_{EDAHFnu} = 3.584$, $p_t=0.063$, $p_{perm}=0.068$). Furthermore, when Player-B is FACE, we find
207 that subjects who rated Player-B high in human-likeness are more likely to experience a stronger emo-
208 tional reaction than participants who rated it low (1.731 vs. -0.129, mean $\text{diff}_{EDAsymp}=-1.859$, $p_t=0.017$,
209 $p_{perm} = 0.000$; 5.865 vs. 0.275 $\text{EDAHFnu}=-5.590$, $p_t=0.009$, $p_{perm} = 0.000$). We do not find a similar effect
210 when Player-B is Human or Computer-box. Finally, we note that the psychophysiological reaction of
211 subjects rating FACE high in human-likeness is significantly higher than that experienced by subjects
212 interacting either with Computer-box or Human, regardless of the rating of human-likeness.

213 Regarding the relationship between the psychophysiological reaction of participants and their choices,
214 we do not find any significant correlation using the two indexes EDAsymp and EDAHFnu. However,
215 if we split our participants into two groups according to whether they express a stronger (or weaker)
216 psychophysiological reaction than the median level of the distribution of EDAsymp (see Table 6), we
217 can observe that those who experienced a stronger reaction are also less likely to choose 'In' in both
218 Computer (0.636 vs. 0.909, mean $\text{diff}=0.273$, and $p_p=0.015$) and Human (0.462 vs. 0.750, $\text{diff}=0.288$, and
219 $p_p=0.070$).

220 Finally to study the interaction between human-likeness and psychophysiological reaction of our
221 participants we conduct a probit analysis for the probability of playing 'In' using as a set of regressors
222 player human-likeness and EDAsymp dummy, along with a dummy for each treatment. Results are
223 reported in Figure 3. This figure highlights that increasing the psychophysiological reaction from low
224 to high reduces the probability of playing 'In'. However, increasing the level of human-likeness coun-
225 terbalances this negative effect, especially in FACE and in Computer-box.

226

227 3 Discussion and conclusion

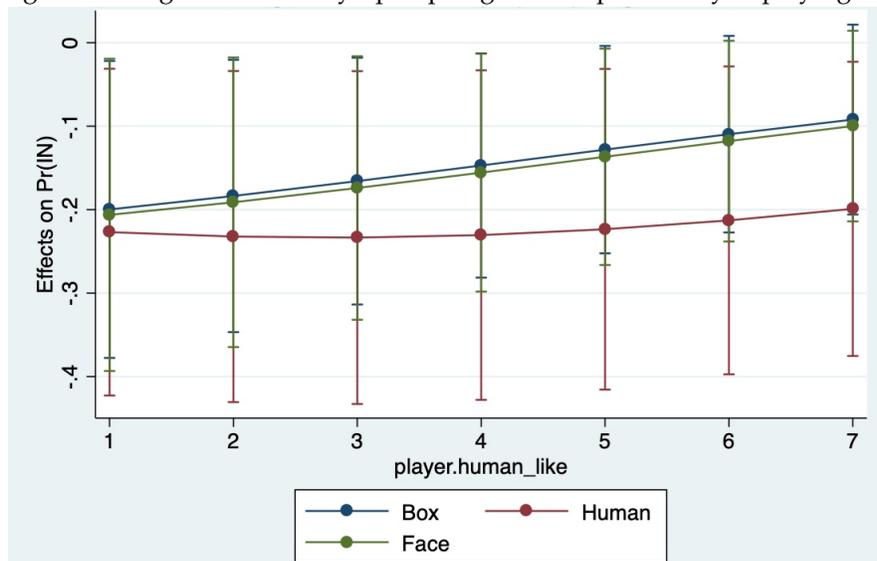
228 In our experiment participants were confronted with a counterpart which differed in the degree of
229 human-likeness: a light-emitting computer-box, a female humanoid and a female human (which re-
230 sembled the humanoid). The participants needed to decide - after listening to a message from the

Table 6: RELATIVE FREQUENCIES OF 'CHOICE IN' BY PHYSIOLOGICAL STATE AND HUMAN-LIKENESS

| | | EDAsymp | | Total |
|--------------|-------|---------------|---------------|---------------|
| | | High | Low | |
| FACE | High | 0.916 [12] | 0.900 [10] | 0.909 [22] |
| | Low | 0.667 [12] | 0.714 [14] | 0.692 [26] |
| | Total | 0.792 [24] | 0.792 [24] | 0.792 [48] |
| Computer-Box | High | 0.667 [7] | 0.857 [9] | 0.750 [16] |
| | Low | 0.616 [15] | 0.933 [13] | 0.786 [28] |
| | Total | 0.636 [22] | 0.909 [22] | 0.770 [44] |
| Human | High | 0.500 [8] | 0.875 [8] | 0.686 [16] |
| | Low | 0.400 [5] | 0.500 [4] | 0.444 [9] |
| | Total | 0.462 [13] | 0.750 [12] | 0.600 [25] |

Each cell represents the frequencies of choice 'In' within each category. An individual is classified in EDAsymp High whenever is above the median level of the EDAsymp distribution, and EDAsymp Low otherwise. Human-likeness is Low when the participant rating is in the lower side of the distribution on the 7-likert scale, and High otherwise. The number of observations are in squared brackets.

Figure 3: Marginal effect of Sympamp High on the probability of playing 'In'



231 counterpart, containing in half of the cases a promise - whether to trust or not their opponent in the
232 game. We find evidence that a human receiving a promise from a humanoid has more trust in it only
233 when he (or she) perceived the artificial agent very similar to a human-being. Indeed, if we replace the
234 social robot with a human we find a similar pattern. However, replacing it by the computer-box the
235 effect of receiving a promise disappears. We also find that participants experienced a stronger psycho-
236 physiological reaction when confronted with a humanoid, especially if it appeared to them very close
237 to human. Moreover, we observe that those participants expressing stronger psychophysiological reac-
238 tion were less likely to trust the counterpart (i.e. chose more often the safer option) when this is either
239 a computer-box or a human.

240 Taken all together, these results suggest that human-likeness and (integral) emotions play both an
241 important role in the decision to trust the counterpart, possibly in interaction with each other. However,
242 some remarks follow in order. While in this experiment we can fully control the degree of human-
243 likeness by varying it across treatments, we have less control over the type of emotions experienced
244 by our subjects. Although physiological measures such electrodermal activity (EDA) have been widely
245 used over the last decades for representing emotional arousal, and most scholars accept a physiological
246 component in the definition of emotions, it is not possible to directly match the physiological state of
247 a participant with a direct type of emotion (e.g. fear or anxiety). In addition, the literature on emotion
248 arousal highlights that there might be individuals exhibiting different physiological responses to the
249 same emotional state [37]. Therefore, our results can only suggest a greater or a weaker 'emotional
250 arousal' without giving any insights into the type of emotions proved by our participants.

251 Nevertheless, the vast psychological literature on emotions and decision-making offers us an inter-
252 esting framework to interpret our results. In particular, recent evidence from laboratory experiments is
253 mostly consistent with the Appraisal-Tendency Framework according to which emotions change indi-
254 viduals' appraisal of a situation, thereby affecting individual choices [9, 38]. Importantly, in that fram-
255 ing, emotions of the same valence (such as fear and anger) can exert opposing influences on choices.
256 Thus, what matters is whether an emotion (either positive or negative) by leading to a more cautious
257 appraisal of the situation reduces the feeling of control, e.g. thereby reducing the willingness to take
258 risks. Therefore, even if we are not able to disentangle among different types of emotions, we can rea-
259 sonably assert that in our framework, whenever the experience of a stronger emotional arousal lead a
260 participant to a more cautious appraisal of the counterpart, we observe a more careful assessment of
261 the situation and a lower willingness to take risk and trust the counterpart. This interpretation of our
262 results is also consistent with previous research showing that participants with ventromedial prefrontal

263 cortex (a key area of the brain for integrating emotion and cognition) repeatedly select a riskier financial
264 option over a safer one, even to the point of bankruptcy, despite their understanding of the suboptim-
265 ality of their choices. In particular, their physiological measure of skin response suggests that they did
266 not experience the emotional signals (i.e. the somatic markers) that lead normal decision makers to
267 fear high risks [9]. Overall, these results strongly support the efforts in developing technologies enhan-
268 cing the humanity of social robots, both in terms of human appearance and communication behaviour.
269 Indeed, if from one-side it is not possible to control for human emotions, in line with recent studies
270 [21, 22], our results suggest that increasing the human-likeness of an artificial agent increases sensibly
271 the likelihood that a human counterpart will trust it. At the same time, the analysis we conducted opens
272 an interesting question about the role of specific emotions, also over the longer time-horizons, that we
273 are not able to fully disentangle in our simple one shot-game.

274 To conclude, we see several directions for future interdisciplinary research. The first one is to explore
275 different types of human-robot interactions, for example, prisoner dilemma games, coordination games
276 or repeated interactions (e.g. by replicating the analysis of Crandall and co-authors with a social robot
277 [19]). The second direction of research is on the side of the social robot. It would be very interesting
278 to introduce - within standard experiments in economics - the behavior of people interacting with a
279 robot that can also additionally adapt its facial expression, as well as the mode of communication, to
280 the perceived emotions of the human counterpart.

281 References

- 282 [1] Lange, P. A. M. V. Generalized trust: Four lessons from genetics and culture. *Current Directions in*
283 *Psychological Science* **24**, 71–76 (2015). (document)
- 284 [2] Fehr, E. On the economics and biology of trust. *Journal of the european economic association* **7**, 235–266
285 (2009). (document)
- 286 [3] Langevoort, D. C. Selling hope, selling risk: some lessons for law from behavioral economics about
287 stockbrokers and sophisticated customers. *Cal L. Rev.* **84**, 627 (1996). (document)
- 288 [4] Nishio, S., Ogawa, K., Kanakogi, Y., Itakura, S. & Ishiguro, H. Do robot appearance and speech
289 affect people,Äôs attitude? evaluation through the ultimatum game. *Geminoid Studies: Science and*
290 *Technologies for Humanlike Teleoperated Androids* 263–277 (2018). (document)
- 291 [5] Picard, R. W. Toward machines with emotional intelligence. In *ICINCO (Invited Speakers)*, 29–30
292 (Citeseer, 2004). (document)
- 293 [6] Engelmann, J. B., Meyer, F., Ruff, C. C. & Fehr, E. The neural circuitry of emotion-induced distor-
294 tions of trust. *BioRxiv* 129130 (2018). (document)
- 295 [7] Schniter, E., Shields, T. W. & Sznycer, D. Trust in humans and robots: Economically similar but
296 emotionally different (2018). (document)
- 297 [8] Jung, E.-S., Dong, S.-Y. & Lee, S.-Y. Neural correlates of variations in human trust in human-like
298 machines during non-reciprocal interactions. *Scientific reports* **9**, 1–10 (2019). (document)
- 299 [9] Lerner, J. S., Li, Y., Valdesolo, P. & Kassam, K. S. Emotion and decision making. *Annual review of*
300 *psychology* **66** (2015). (document), 3
- 301 [10] Damasio, A. R. The somatic marker hypothesis and the possible functions of the prefrontal cortex.
302 *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences* **351**, 1413–1420
303 (1996). (document)
- 304 [11] Damasio, A. R. Descartes error revisited. *Journal of the History of the Neurosciences* **10**, 192–194
305 (2001). (document)
- 306 [12] Vaa, T. Driver behavior models and monitoring of risk: Damasio and the role of emotions. In *In-*
307 *ternational Conference: Traffic Safety on Three Continents*PTRC Education and Research Services Limited,
308 VTI Konferens 18A (2001). (document)

309 [13] Fox, A. S., Lapate, R. C., Shackman, A. J. & Davidson, R. J. *The nature of emotion: fundamental*
310 *questions* (Oxford University Press, 2018). (document)

311 [14] Arkin, R. C., Ulam, P. & Wagner, A. R. Moral decision making in autonomous systems: Enforce-
312 ment, moral emotions, dignity, trust, and deception. *Proceedings of the IEEE* **100**, 571–589 (2011).
313 (document)

314 [15] Tortosa, M. I., Strizhko, T., Capizzi, M. & Ruz, M. Interpersonal effects of emotion in a multi-round
315 trust game. *Psicologica: International Journal of Methodology and Experimental Psychology* **34**, 179–198
316 (2013). (document)

317 [16] Campellone, T. R. & Kring, A. M. Who do you trust? the impact of facial emotion and behaviour
318 on decision making. *Cognition & emotion* **27**, 603–620 (2013). (document)

319 [17] Engelmann, J. B. *et al.* Emotions can bias decision-making processes by promoting specific behavi-
320 oral tendencies (2018). (document)

321 [18] Hancock, P. A. *et al.* A meta-analysis of factors affecting trust in human-robot interaction. *Human*
322 *factors* **53**, 517–527 (2011). (document)

323 [19] Crandall, J. W. *et al.* Cooperating with machines. *Nature communications* **9**, 1–12 (2018). (document),
324 3

325 [20] Nitsch, V. & Glassen, T. Investigating the effects of robot behavior and attitude towards technology
326 on social human-robot interactions. In *2015 24th IEEE International Symposium on Robot and Human*
327 *Interactive Communication (RO-MAN)*, 535–540 (IEEE, 2015). (document), 1

328 [21] Waytz, A., Heafner, J. & Epley, N. The mind in the machine: Anthropomorphism increases trust in
329 an autonomous vehicle. *Journal of Experimental Social Psychology* **52**, 113–117 (2014). (document), 3

330 [22] Nass, C. & Moon, Y. Machines and mindlessness: Social responses to computers. *Journal of social*
331 *issues* **56**, 81–103 (2000). (document), 3

332 [23] March, C. The behavioral economics of artificial intelligence: Lessons from experiments with
333 computer players (2019). (document)

334 [24] Mori, M. The uncanny valley: The original essay by masahiro mori. *IEEE Robots &* (2017). (docu-
335 ment)

- 336 [25] Charness, G. & Dufwenberg, M. Promises and partnership. *Econometrica* **74**, 1579–1601 (2006).
337 (document), 1
- 338 [26] Berg, J., Dickhaut, J. & McCabe, K. Trust, reciprocity, and social history. *Games and economic behavior*
339 **10**, 122–142 (1995). (document)
- 340 [27] Tao, J. & Tan, T. Affective computing: A review. In *International Conference on Affective computing*
341 *and intelligent interaction*, 981–995 (Springer, 2005). (document)
- 342 [28] Mazzei, D. *et al.* The face of autism. In *19th International Symposium in Robot and Human Interactive*
343 *Communication*, 791–796 (IEEE, 2010). 1
- 344 [29] Lazzeri, N. *et al.* Can a humanoid face be expressive? a psychophysiological investigation. *Frontiers*
345 *in bioengineering and biotechnology* **3**, 64 (2015). 1
- 346 [30] Cominelli, L., Mazzei, D. & De Rossi, D. E. Seai: Social emotional artificial intelligence based on
347 damasio,Äôs theory of mind. *Frontiers in Robotics and AI* **5**, 6 (2018). 1, 4.3
- 348 [31] Posada-Quintero, H. F. *et al.* Power spectral density analysis of electrodermal activity for sympath-
349 etic function assessment. *Annals of biomedical engineering* **44**, 3124–3135 (2016). 1, 4.6
- 350 [32] Acharya, U. R., Joseph, K. P., Kannathal, N., Lim, C. M. & Suri, J. S. Heart rate variability: a review.
351 *Medical and biological engineering and computing* **44**, 1031–1051 (2006). 1, 4.6
- 352 [33] Ghiasi, S. *et al.* A new sympathovagal balance index from electrodermal activity and instantaneous
353 vagal dynamics: A preliminary cold pressor study. In *2018 40th Annual International Conference of*
354 *the IEEE Engineering in Medicine and Biology Society (EMBC)*, 3068–3071 (IEEE, 2018). 1, 4.7
- 355 [34] Ghiasi, S., Greco, A., Barbieri, R., Scilingo, E. P. & Valenza, G. Assessing autonomic function from
356 electrodermal activity and heart rate variability during cold-pressor test and emotional challenge.
357 *Scientific Reports* **10**, 1–13 (2020). 1
- 358 [35] Franke, T., Attig, C. & Wessel, D. A personal resource for technology interaction: development
359 and validation of the affinity for technology interaction (ati) scale. *International Journal of Human–*
360 *Computer Interaction* **35**, 456–467 (2019). 1, 3
- 361 [36] Rieger, M. O., Wang, M. & Hens, T. Risk preferences around the world. *Management Science* **61**,
362 637–648 (2015). 1

- 363 [37] Picard, R. W. *Affective computing* (MIT press, 2000). 3
- 364 [38] Meier, A. N. Emotions, risk attitudes, and patience. Tech. Rep., SOEPpapers on Multidisciplinary
365 Panel Data Research (2019). 3
- 366 [39] Greiner, B. *et al.* The online recruitment system orsee 2.0-a guide for the organization of experi-
367 ments in economics. *University of Cologne, Working paper series in economics* **10**, 63–104 (2004). 4.1
- 368 [40] Lazzeri, N., Mazzei, D., Cominelli, L., Cisternino, A. & De Rossi, D. E. Designing the mind of a
369 social robot. *Applied Sciences* **8**, 302 (2018). 4.3
- 370 [41] Bosse, T., Jonker, C. M. & Treur, J. Formalisation of damasio,Äôs theory of emotion, feeling and
371 core consciousness. *Consciousness and cognition* **17**, 94–113 (2008). 4.3
- 372 [42] Zarakı, A. *et al.* Design and evaluation of a unique social perception system for human–robot
373 interaction. *IEEE Transactions on Cognitive and Developmental Systems* **9**, 341–355 (2016). 4.3
- 374 [43] Cominelli, L. *et al.* A preliminary framework for a social robot ,Äúsixth sense,Äù. In *Conference on*
375 *Biomimetic and Biohybrid Systems*, 58–70 (Springer, 2016). 4.3
- 376 [44] Mazzei, D., Cominelli, L., Lazzeri, N., Zarakı, A. & De Rossi, D. I-clips brain: A hybrid cognit-
377 ive system for social robots. In *Conference on Biomimetic and Biohybrid Systems*, 213–224 (Springer,
378 Cham, 2014). 4.3
- 379 [45] Giarratano, J. C. & Riley, G. *Expert systems* (PWS publishing co., 1998). 4.3
- 380 [46] Russell, J. A. A circumplex model of affect. *Journal of personality and social psychology* **39**, 1161 (1980).
381 4.3
- 382 [47] Cominelli, L. *et al.* Damasio,Äôs somatic marker for social robotics: preliminary implementation
383 and test. In *Conference on Biomimetic and Biohybrid Systems*, 316–328 (Springer, 2015). 4.3
- 384 [48] Mazzei, D., Lazzeri, N., Hanson, D. & De Rossi, D. Hefes: An hybrid engine for facial expressions
385 synthesis to control human-like androids and avatars. In *2012 4th IEEE RAS & EMBS International*
386 *Conference on biomedical robotics and biomechatronics (BioRob)*, 195–200 (IEEE, 2012). 4.3
- 387 [49] Kreibig, S. D. Autonomic nervous system activity in emotion: A review. *Biological psychology* **84**,
388 394–421 (2010). 4.6

- 389 [50] Vernet-Maury, E., Deschaumes-Molinari, C., Delhomme, G. & Dittmar, A. Autonomic nervous
390 system activity and mental workload. *International Journal of Psychophysiology* **14**, 153–154 (1993).
391 4.6
- 392 [51] Greco, A., Valenza, G., Bicchi, A., Bianchi, M. & Scilingo, E. P. Assessment of muscle fatigue during
393 isometric contraction using autonomic nervous system correlates. *Biomedical Signal Processing and*
394 *Control* **51**, 42–49 (2019). 4.6
- 395 [52] Greco, A., Valenza, G. & Scilingo, E. P. *Advances in Electrodermal activity processing with applications*
396 *for mental health* (Springer, 2016). 4.6
- 397 [53] Greco, A., Valenza, G., Lanata, A., Scilingo, E. P. & Citi, L. cvxeda: A convex optimization approach
398 to electrodermal activity processing. *IEEE Transactions on Biomedical Engineering* **63**, 797–804 (2015).
399 4.6
- 400 [54] Strigo, I. A. & Craig, A. D. Interoception, homeostatic emotions and sympathovagal balance.
401 *Philosophical Transactions of the Royal Society B: Biological Sciences* **371**, 20160010 (2016). 4.7
- 402 [55] Sleight, P. & Bernardi, L. Sympathovagal balance. *Circulation* **98**, 2640–2640 (1998). 4.7

403 **4 Methods**

404 **4.1 Participants**

405 The experimental protocol was approved with unanimity by the Bioethical Committee of the University
406 of Pisa (Review No. 21/2019), and all sessions were performed in accordance with relevant regulations
407 and guidelines. Informed consent was obtained from all participants in the experiments. Informed
408 consent was also obtained from the actress to publish online Figure 1.

409 Participants were invited through ORSEE system of the University of Pisa, which allow to randomly
410 invite participants and to keep track of their participation in experiments [39]. In total 164 participants
411 signed-up and showed up in the laboratory in the day they were invited. Two subjects were removed
412 from the pool because they did not followed the experimental procedure correctly. The final sample
413 was therefore of 162 (90 students were female and 72 male, with a mean age of about 26 years old).

414 The total number of participants was determined based on previous study (in Charness and Dufwen-
415 berg 2006) as well as on technical constraints (i.e. the possibility to run the humanoid for a long period
416 of time). In Charness and Dufwenberg (2006) there were 42 pairs in the session in which participants

417 could receive a message from the opponent, with a share of 0.74 of player-B choosing 'In'. We knew that
418 given the proportion of 0.74 in the human treatment, the smallest difference that could be detected with
419 this sample size and a power of 0.80 was about 0.20. Therefore, we aim to have a final sample of about
420 50 participants, thereby we aimed to recruit 55 participants for each treatment. In the computer-box
421 machine we recruited slightly more subjects as we did experience several participants' no-show up in
422 the previous treatment.

423 **4.2 Experimental procedure**

424 Each participant arrives in the laboratory and enter a room in which (s)he read and sign the consent
425 to participate in the study. The participant then sits in front of a computer screen where (s)he can read
426 autonomously the experiment instructions and fill in some preliminary questions about their attitudes
427 towards the technology. At this stage participant also worn the bracelet on their left wrist. This phase
428 will be then used as 'the rest' phase for measuring psychophysiological parameter (see also below).
429 Once the time dedicated to this part has expired, the participant is lead by the experimenter to another
430 room where the robot is located and a marker is recorded on the bracelet to begin the second phase of
431 measurement of psychophysiological parameters. The participant seats on chair, always located at the
432 same distance from the robot, and when is ready to start the experiment has to rise his hand. At this
433 point, the robot welcomes the participant with a standard sentence ('Nice to meet you! Let's start') to
434 then state one random sentence out of 8 (according to the treatment, see again Table 1 in the paper).
435 The robot then tells the participant a standard final sentence, inviting the participant to enter his(her)
436 choice in the computer in front of him(her). The robot cannot observe though the choice the participant
437 has made. To conclude the experiment, the participant has to return to the initial room, to complete an
438 exit questionnaire about the interaction of the robot, and receive the final payment.

439 **4.3 The FACE Robot and the SEAI Cognitive System**

440 The FACE robot (Facial Automaton for Conveying Emotions) is a humanoid with hyper-realistic adult
441 female aesthetics, specifically designed for social robotics [40]. It is composed with a passive body
442 on the top of which a Hanson Robotics' head has been mounted. The head is designed to host 32
443 servomotors that guide the neck of the robot, its eyes, mouth, and facial expression. The face of the
444 ginoid is made of Frubber,³ a registered material with skin-like mechanical and aesthetical features.

³<https://patents.google.com/patent/US7113848?q=frubber>

445 This hardware is controlled by SEAI (Social Emotional Artificial Intelligence), a distributed control ar-
446 chitecture made of perception, cognitive and actuation systems, that endow the robot with expressive
447 and communicative capabilities [30], including also the possibility to emulate verbal communication
448 following prerecorded audio files.⁴ SEAI is a bio-inspired architecture based on neuroscientific theories
449 of mind. In particular, it has been inspired by the findings of Antonio Damasio and it is consistent
450 with the computational formalization made by [41]. In its development, the influence of emotions in
451 the decision-making process has been of primary importance. The perception part of the system is the
452 Scene Analyzer, an audiovisual perception system conceived to analyze a social environment using the
453 robot sensors and to extract meaningful social cues from these available data. Features that can be ex-
454 tracted from a human interlocutor are, e.g., the three dimensional position of 25 joint coordinates, their
455 speaking probability, meaningful postures and gestures, estimated facial expressions, age and gender
456 [42]. This Social Perception System has already been successfully integrated with the acquisition of
457 physiological parameters (i.e., electrodermal activity, respiration rate and heart rate variability) in past
458 experiments [43]. All the environmental information analyzed by the perception system of the robot is
459 then processed by its cognitive system, i.e., the I-CLIPS Brain [44], a rule-based expert system written in
460 CLIPS language [45]. The knowledge base of the expert system is written by means of IF-THIS-THEN-
461 THAT rules, where each rule contains a set of actions that will be executed if several conditions about
462 the upcoming factual information are satisfied. Thanks to these rules it is possible to design the beha-
463 vior of the humanoid. For example, a particular expression gathered in its interlocutor can lead to the
464 trigger of a sentence or a facial expression performed by the robot, but also to the modification of the ro-
465 bot's internal values. In fact, SEAI includes emotional internal values (i.e., valence and arousal), which
466 combination describes an emotional state, here defined as *mood*. This method of representing emotion
467 is based on the well-known Russell's Circumplex Model of Affect [46]. In the case of the robot, mood
468 is not necessarily externalised by perceivable movements, rather it is implied in biasing the chaining
469 of the rules, and so, the decision tree of the robot. Emotion biasing decision in this cognitive system
470 has been previously tested [47]. The instructions coming from the cognitive block about the emotion to
471 be expressed through facial expression - (v,a) values, the sentence to say, and the point to look at, are
472 merged and continuously executed thanks to the actuation system, which translate them in movements
473 performed by the motors that drive the face, the mouth and the neck of the humanoid [48]. Furthermore,
474 the SEAI architecture is completely modular and portable, all the blocks composing the framework are

⁴The audio files used for the experiment have been recorded using the voice of a professional actress, the same who interpreted the role of Player-B in the interactions with the real person; the sentences were the Italian translation of the sentences between the Charness trust game players.

475 stand-alone applications that process a limited set of information. These modules are distributed in a
476 local net of computers that communicate by means of the YARP middleware.⁵ This implies that each
477 module can be activated or deactivated, and that the perception and cognitive systems can be used also
478 without controlling the FACE Robot. As a result, we were able to use exactly the same rules engine
479 in the computer box case, simply disabling the actuation part of the system that control the robot, and
480 using instead the bluetooth speaker, presented as a smart computer box, actually running the same
481 perception and actuation system of the robot. This led to a very close and controlled comparison.

482 **4.4 How the robot takes a decision, the Rules Engine**

483 In this experiment, the robot (as well as the computer box) decides whether to *Roll* or *Don't Roll* accord-
484 ing to its emotional state and following its decision rules. In particular, a positive mood in SEAI (i.e., an
485 emotional state with positive valence) will lead the robot to be collaborative with the human player and
486 play *Roll*; while a negative mood in SEAI (i.e., an emotional state with negative valence) will lead the
487 robot to play *Don't Roll* (see Figure 5). The decision is taken at the end of the interaction with Player-A,
488 when the subject goes out of the room, and so out of the field of view of the robot.

489 If in the moment in which the robot has to take a decision, it is in a qualitatively neutral mood ($v=0$,
490 regardless the arousal), the decision will be taken randomly (50%). Participants' behavior during all
491 the time spent alone in the room with the robot, once observed by the Scene Analyzer and processed
492 in SEAI, act as an input modifying the parameters of the robot which correspond to its 'mood', thus in
493 turn affecting its course of action (i.e., its final decision). However, in this experiment, at each interac-
494 tion with a new participant the robot always resetted its internal values at the «neutral emotional state»
495 (which corresponds to $v = 0, a = 0$ in the graph). In conclusion, thanks to SEAI the robot was com-
496 pletely autonomous, by means of the rules everything was pre-programmed and automatized, starting
497 from the rules that use perceived social cues to modulate the emotional state of the robot, to other rules
498 determining which sentence it has to say, when to start and to end a treatment, and the storage of all
499 the data acquired with timestamps in a structured dataset. The complete code of the rules engine is
500 available upon request from the authors.

⁵<https://www.yarp.it/>

501 4.5 Mean comparisons across groups

502 To compare the means (μ) of the distribution of a random variable for two independent groups (X, Y),
503 we perform *t-Student* tests on the equality of means. Specifically, to test for $\mu_x = \mu_y$ (when the variances
504 σ_x and σ_y are unknown and replaced by s_x and s_y) the test is $t = \frac{\bar{x} - \bar{y}}{(\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y})^{1/2}}$ which is distributed as Stu-
505 dent's t . When the random variable is not continuous but a proportion, we use a normally distributed
506 test statistic calculated as $z = \frac{\hat{p}_x - \hat{p}_y}{(\hat{p}_q(1-\hat{p}_q)(1/n_1 + 1/n_2))^{1/2}}$ where $p_x = \frac{x+y}{n_1+n_2}$ where x and y are the number
507 of successes in the two populations.

508 Both t and proportion tests rely on assumption about the distribution of the data. This is the reason
509 why we also rely on permutation tests, which are nonparametric tests - i.e. do not rely on any assump-
510 tion about the distribution of the data. Permutation tests work by resampling the observed data many
511 times. The permutation test based on means implies: 1) to compute the sample means for each group
512 $d_{observed} = \bar{x} - \bar{y}$; 2) pool all the data together and randomly permute the pooled data; 3) then com-
513 pute again the sample mean again for the two groups and note the difference d_1 ; 3) repeat step 2 and
514 3 several times in order to obtain several mean differences, i.e. d_1, d_2, d_3, \dots . If the null hypothesis of no
515 difference between the two groups is true, by changing the order of the data we should not observe any
516 difference in the means, otherwise it should look different from the real data. The ranking of the real
517 test statistic, i.e. $d_{observed}$, among the shuffled test statistics, d_1, d_2, d_3, \dots , gives a p -value.

518 4.6 Description and analysis of Physio data

519 Pulse rate variability (PRV) and electrodermal activity (EDA) signals are directly modulated by the
520 autonomic nervous system (ANS) activity and, therefore, are considered ideal non-invasive physiolo-
521 gical signals to investigate the ANS dynamics. Indeed, the ANS plays a crucial role in the processing of
522 the emotional response, mental fatigue and workload [49, 50, 51].

523 Particularly, the EDA signal measures the activity of eccrine sweat glands on the hand surface.
524 Since sweat glands are directly innervated by the sympathetic branch of the ANS (and in particular
525 the sudomotor nerve), the EDA analysis is considered one of the best ways to monitor the sympathetic
526 activity [52]. As a preprocessing step, we applied the well-known cvxEDA model [53] to remove the
527 superimposed noise. From each free-to-noise EDA signal, we estimated the power spectrum within
528 the frequency range of 0.045 and 0.25Hz (EDAsymp), which has been demonstrated to be an effective
529 estimator of the sympathetic nervous system activity [31].

530 The PRV signal was computed interpolating the interbeat interval time series (IBI) extracted from the

531 photoplethysmography signals acquired by the Empatica wearable acquisition system. To characterize
532 the activity of the parasympathetic nervous system, which, as known, regulates the high frequency
533 oscillations of the PRV signal, we estimated the Power Spectral Density (PSD) related to each PRV
534 signal [32]. Two main spectral bands were considered: low frequency (LF) band (ranging between 0.04
535 and 0.15 Hz), and high frequency (HF) band (from 0.15 to 0.4 Hz). Then, the power spectrum in the HF
536 band normalized to the sum of LF and HF power (HFnu) was computed to quantify the activity of the
537 parasympathetic nervous system.

538 Note that all physiological indexes computed during the interaction with the agent were normalized
539 for each participant by dividing them by the baseline value computed before the interaction phase

540 **4.7 New index from the sympathovagal assessment**

541 Emotions regulation process modulates the sympathovagal balance [54, 55], which is considered a re-
542 liable marker of the human affective state. Previous studies have suggested that LF power spectrum
543 can provide a quantitative marker of the sympathetic outflow and have used the LF/HF ratio as a cor-
544 relate of the sympathovagal balance. However, the LF power is now regarded as a measure of both
545 sympathetic and vagal tone, leading to ambiguities and possible inconsistent conclusions on the use of
546 the LF/HF ratio as sympathovagal marker. In this study, we employed novel indexes of the sympath-
547 ovagal dynamics based on the combination of the information extracted from the EDA and PRV signal
548 [33]. Indeed, while EDAsymp reliably characterizes the sympathetic activity, HFnu is considered an
549 effective cardiovascular-related features it that reliably quantify the parasympathetic outflow. Accord-
550 ingly, we have estimated the sympathovagal balance using the ratio between EDAsymp and HFnu:
551 EDAsymp/HFnu [33].

Figure 4: EMOTIONAL STATE OF THE ROBOT

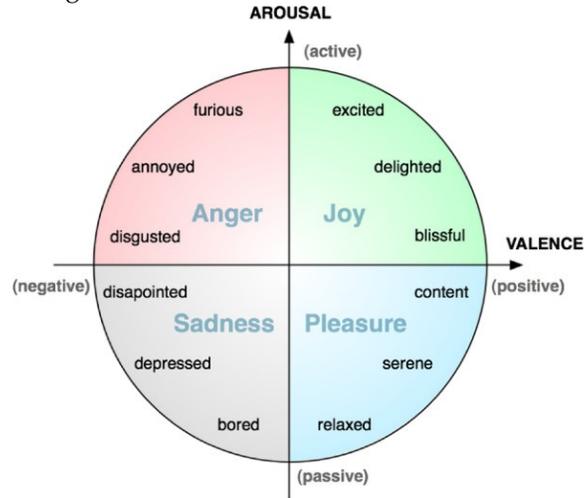
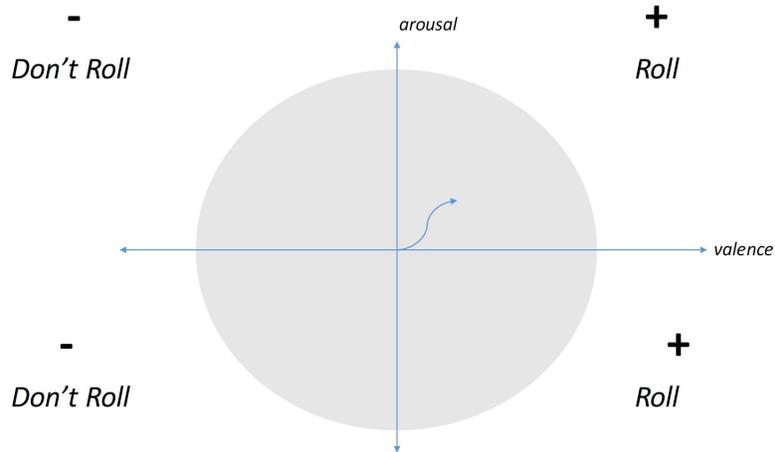


Figure 5: DECISION RULE OF THE ROBOT



INSTRUCTIONS: English translation from Italian

Welcome! This experiment will last about 30 minutes. You will receive 5 Euro for your participation. Based upon the choices you will take in the experiment; you can earn additional money. We now ask you to turn off your mobile phone and to read the instructions carefully.

The aim of this experiment is to study how people take decisions. In particular, this experiment wants to study how people take decision when interacting with a human-like robot.

Should you have any doubt, please do not hesitate to ask clarifications to the experimenter.

The data related to this experiment will be saved and analyzed anonymously. No video will be recorded.

In this experiment you will play with FACE i.e. a social robot which is able to prove and express its emotions. [with a computer-box which is given a system of social perception]. FACE [The Computer box] is also able to take its decisions autonomously, following its own behavioral rules. In this game, FACE [The Computer box] is programmed to choose autonomously between two actions: **ROLL** and **DON'T ROLL** a six-faces dice.

[In this experiment you will play with Deborah. Deborah can choose autonomously between two actions: **ROLL** and **DON'T ROLL** a six-faces dice.]

YOUR CHOICE

You will have to choose between two options: whether to play **IN** or **OUT**.

Should you choose **OUT**, both you and FACE [Computer box] [Deborah] will earn 5 Euro each.

Should you choose **IN**, FACE [Computer box] [Deborah] can then choose between the two options: **ROLL** and **DON'T ROLL** the six-faces dice. In the event FACE [Computer box] [Deborah] choosing **DON'T ROLL**, you will receive 0 Euro and FACE [Computer box] [Deborah] will earn 14 Euro. In the event FACE [Computer box] [Deborah] choosing **ROLL**, FACE [Computer box] [Deborah] will always earn 10 Euro

while your earning depends on the results of the dice roll. If the result of the dice roll is a number between 2 and 6 you will earn 12 Euro, otherwise if the result of the dice roll is the number 1 you will receive 0 Euro.

It is important to notice that FACE [Computer box] [Deborah] will not know whether you opted either IN or OUT when it has to reach a decision. It is also important to notice that the money earned by FACE will remain to FACE itself [will remain to the lab (e.g. maintenance)], and used for its necessity (e.g. maintenance)

The payments are summarized in the table below.

| | Dice roll | You earning | FACE's [Computer box] [Deborah] earning |
|--|-----------------------|-------------|---|
| If you choose OUT | - | 5Euro | 5Euro |
| If you choose IN FACE choose DON'T ROLL | - | 0Euro | 14 Euro |
| If you choose IN FACE choose ROLL | Result: 1 | 0Euro | 10 Euro |
| | Results: 2,3,4,5,6 | 12 Euro | 10 Euro |

Now you have 5 minutes to read these instructions alone and ask clarifications questions to the experimenter. Once you have finished reading, the experimenter will bring you to another room where FACE [Computer box] [Deborah] is. You will have to sit on the chair in front of FACE, and in order to begin the experiment you need to raise your right hand. At the point, you will hear a message from FACE [Computer box] [Deborah]. You will then enter your choice in the computer close to you.

Once you have done, we will wait for you to come back again to this room, to fill in a final questionnaire and receive your payment.

Figures



Figure 1

THREE TYPES OF PLAYER-B

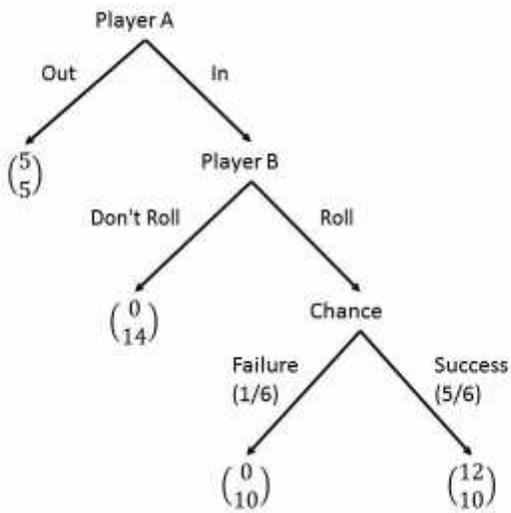


Figure 2

THE GAME

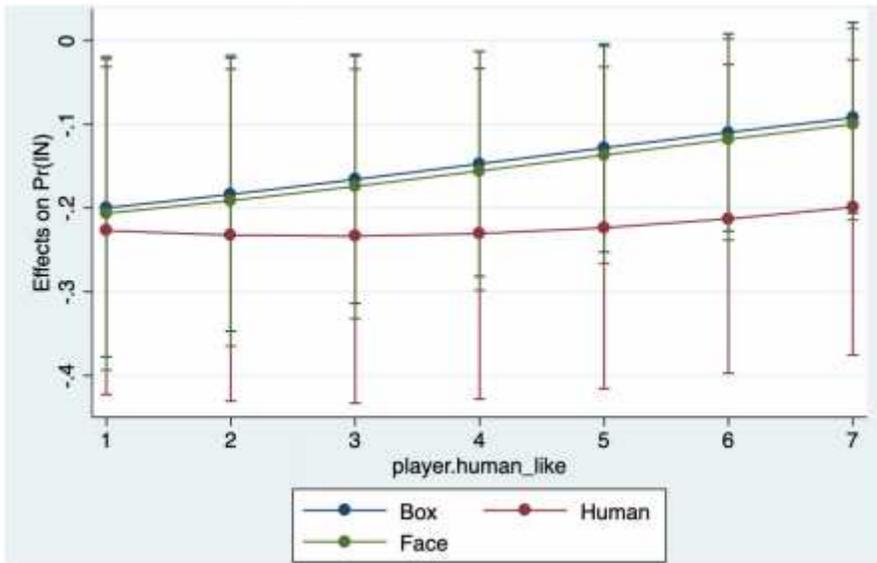


Figure 3

Marginal effect of Sympamp High on the probability of playing 'In'

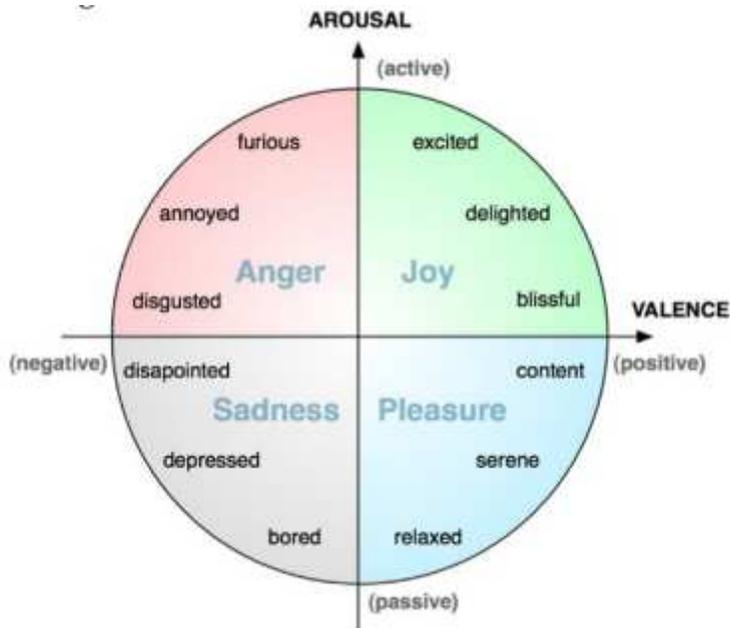


Figure 4

EMOTIONAL STATE OF THE ROBOT

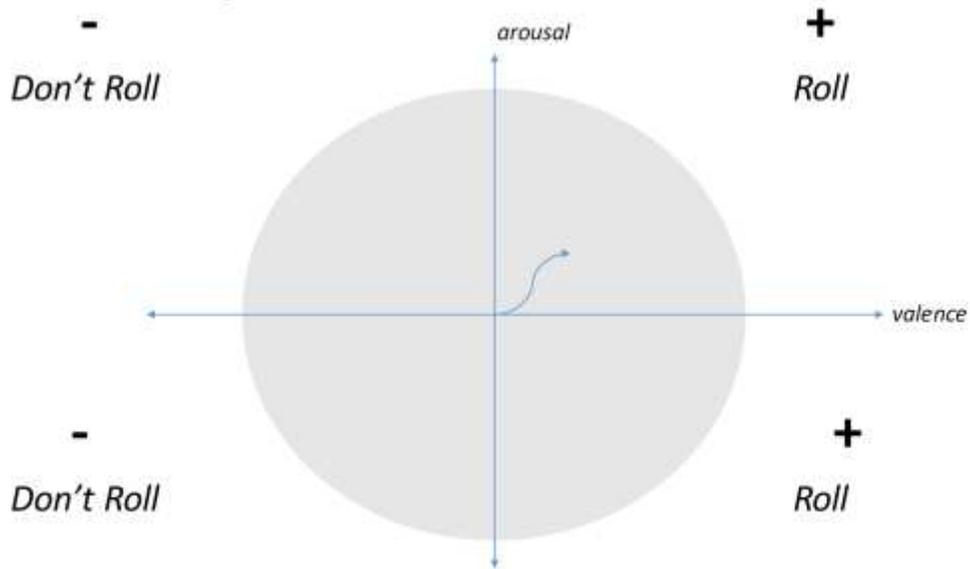


Figure 5

DECISION RULE OF THE ROBOT