

Personality Perception of Robot Avatar Tele-operators in Solo and Dyadic Tasks

Paul Bremner^{1,*†}, Oya Celiktutan^{2,*}, and Hatice Gunes²

¹*Bristol Robotics Laboratory, University of West England, Bristol, UK*

²*Computer Laboratory, University of Cambridge, Cambridge, UK*

Correspondence[†]:

Paul Bremner

Bristol Robotics Laboratory, University of West England, Bristol, United Kingdom,

Paul.Bremner@brl.ac.uk

2 ABSTRACT

3 Humanoid robot avatars are a potential new tele-communication tool whereby a user is remotely
4 represented by a robot that replicates their arm, head and possibly face movements. They have
5 been shown to have a number of benefits over more traditional media such as phones or video
6 calls. However using a tele-operated humanoid as a communication medium inherently changes
7 the appearance of the operator, and appearance based stereotypes are used in interpersonal
8 judgements (whether consciously or unconsciously). One such judgement that plays a key role
9 in how people interact is personality. Hence, we have been motivated to investigate if and how
10 using a robot avatar alters the perceived personality of tele-operators. To do so we carried out
11 two studies where participants performed 3 communication tasks, solo in study one and dyadic in
12 study two, and were recorded on video both with and without robot mediation. Judges recruited
13 using online crowdsourcing services then made personality judgements of the participants in the
14 video clips. We observed that judges were able to make internally consistent trait judgements
15 in both communication conditions. However, judge agreement was affected by robot mediation,
16 although which traits were affected was highly task dependent. Our most important finding was
17 that in dyadic tasks personality trait perception was shifted to incorporate cues relating to the
18 robot's appearance when it was used to communicate. Our findings have important implications
19 for tele-presence robot design and personality expression in autonomous robots.

20

21 **Keywords:** Telepresence, Big Five personality traits, personality perception

1 INTRODUCTION

22 Telecommunication is omnipresent in today's society, with people desiring to be able to communicate
23 with one another, regardless of distance, for a variety of social and practical reasons. While video enabled
24 communication offers a number of benefits over voice only communication, it is still lacking compared to
25 face-to-face interactions Daly-Jones et al. (1998). For example remotely located team members are less
26 included in co-operative activities than co-located team members Daly-Jones et al. (1998), and have fewer
27 conversational turns and speaking time in group conversations O'Conaill et al. (1993). Suggested reasons

*Indicates equal contribution.

28 for these disparities are a lack of social presence of these remote group members, reduced engagement, and
29 reduced awareness of actions Tang et al. (2004). A suggested underlying cause for the disparities found in
30 traditional tele-communication is a lack of physical presence. An alternative is the use of tele-operated
31 robots as communication media. A common approach to such embodied telecommunication is the use of
32 mobile remote presence (MRP) devices: a screen displaying the operators face mounted on a stalk attached
33 to a wheeled base Kristoffersson et al. (2013). Though studies examining the utility of MRPs have found
34 that there are some improvements in social presence, different social norms are observed when people use
35 them to interact, and there are impacts on trust and rapport Rae et al. (2013); Lee and Takayama (2011).
36 Further, such systems are not able to effectively transmit non-verbal communication cues, a key element
37 of human communication not only for information conveyance, but also in maintaining engagement and
38 building rapport Salam et al. (2016).

39 A proposed method for further improving social presence and effectively transmitting body language is to
40 use a humanoid robot as a communication medium. In such a system the operator's body language
41 is duplicated on a humanoid robot such that it is comprehensible and highly salient Bremner and
42 Leonards (2016); Bremner et al. (2016b). Using a humanoid robot as a communications avatar has
43 benefits with regards to engagement of conversational partners Hossen Mamode et al. (2013), social
44 presence Adalgeirsson and Breazeal (2010), group interaction Hossen Mamode et al. (2013), and trust
45 Bevan and Stanton Fraser (2015).

46 However, when using a robot as a remote proxy for communication the operator is represented with a
47 different physical appearance, much as computer generated avatars do in virtual environments. Appearance
48 has been observed to be utilised in making interpersonal judgements Naumann et al. (2009), and this can
49 extend to virtual avatars Wang et al. (2013); Fong and Mar (2015). It was observed that judges made
50 relatively consistent inferences based on avatar appearance alone Wang et al. (2013); Fong and Mar (2015),
51 and more attractive avatars were rated more highly in an interview scenario Behrend et al. (2012). How this
52 might manifest with robot avatars, in particular in the interaction between a robot appearance and human
53 voice communication, remains unclear and is yet to be explored.

54 Here the particular judgement we are concerned with is that of personality perception, an important facet
55 of communication. Researchers in psychology have shown that personality plays a key role in forming
56 interpersonal relationships, and predicting future behaviours Borkeu et al. (2004). These findings have
57 motivated a significant body of work for how people judge others' personalities based on their observable
58 behaviours. A key component of these social cues for personality are non-verbal behaviours. We aim to
59 investigate if such non-verbal personality cues transmitted by a tele-operated humanoid robot continue
60 to be utilised in personality judgements, and how they interact with verbal cues. Non-verbal cues can be
61 transmitted as our robot tele-operation system utilises a motion capture based approach so that arm and
62 head movements the operator performs while talking are recreated with minimal delay on a NAO humanoid
63 robot Bremner and Leonards (2016). The control system is intuitive and immersive, and we observe people
64 behaving similarly to how they do face to face Bremner et al. (2016b).

65 We designed two experiments which follow an experimental methodology common in the personality
66 analysis literature, i.e., videos of participants performing different communication tasks are shown to
67 external observers (judges) for personality assessment (e.g., Borkeu et al. (2004)). Personality judgements
68 are made on the so called big five traits, *extroversion*, *conscientiousness*, *agreeableness*, *neuroticism*,
69 and *openness* (multiple questions relate to each trait). We varied communication media between judges,
70 either video only or robot mediated (also recorded on video). Two main measures are used to see whether
71 there was an effect of communication condition on personality judgements: 1) judge consistency in how

72 they evaluate a given trait, both within and between judge (low consistency indicates lack of cues or
73 conflicting cues); and 2) personality shifts between high and low classification for each trait between the
74 video and robot conditions.

75 Hence we address the following research questions:

- 76 • **RQ1.** Are there differences in judges' consistency in assessing personality traits (within-judge
77 consistency)?
- 78 • **RQ2.** Are there differences in how much judges agree with one another on personality judgements
79 (between-judge consistency)?
- 80 • **RQ3.** Are personality judgements less accurate compared to self ratings (self-other agreement)?
- 81 • **RQ4.** Are perceived personalities systematically shifted to incorporate characteristics associated with
82 the robot's appearance (personality shifts)?

83 This paper is an extended version of our work published in Bremner et al. (2016a). We extended our
84 previous work by adding a second experiment which refined our experimental procedure and used dyadic
85 rather than solo tasks. Our discussions and conclusions are extended to include both experiments, evaluating
86 all our results to give a clearer picture.

87 In the first experiment three tasks are performed direct to camera, i.e., solo tasks. In the second experiment
88 participants performed three tasks that involved interaction with a confederate, i.e., dyadic. The first
89 experiment provided some limited evidence for shifts in personality perception. Further, by adding an audio
90 only communication condition we were able to show that the robot was not simply ignored, and gesture
91 cues performed on the robot were utilised. An important finding from the first experiment was that effects
92 were very task dependent, as the literature suggested. Borkenau et al. (2004) found that *openness* is better
93 inferred in more ability-demanding tasks such as pantomime task. Hence, the second experiment used
94 additional tasks, which by being dyadic will engender personality cues differently; it is also a refinement
95 of our experimental procedure, improving the reliability of our results. It produced compelling evidence
96 that cues related to the robot's appearance were incorporated in personality judgements, causing consistent
97 shifts in perceived personality.

2 RELATED WORK

98 A common approach to investigating personality judgements is first impression or thin slice personality
99 analysis. It is a body of research that studies the accuracy with which people are able to make personality
100 judgements of others based only on short behavioural episodes (termed thin slices). This approach is
101 taken as it is believed that these judgements provide insight into the assessments people make in everyday
102 interactions Funder and Sneed (1993); Borkenau et al. (2004). In such studies, targets are typically asked to
103 perform a range of communication tasks, either solo performances to camera or dyadic with confederates,
104 and are filmed while doing so. *Judges* then observe the video clips and complete personality assessment
105 questionnaires. Ratings of judges are compared with target self ratings, acquaintance ratings, and for
106 inter-judge agreement. For many traits there is sufficient inter-judge agreement for the method to be useful
107 in assessing the impressions a person creates on those they interact with Borkenau et al. (2004); however,
108 the accuracy of judge ratings to self/acquaintance ratings is typically a lot lower, as self/acquaintance
109 ratings are error prone, and use different sources to make their judgements Vinciarelli and Mohammadi
110 (2014).

111 Often analysed in thin slice personality studies are the cues that appear to be utilised in people making
112 their judgements. Appearance, speaking style, gaze, head movements and hand gestures have been
113 frequently reported to be significant predictors of personality Riggio and Friedman (1986); Borkenau et al.
114 (2004); Borkenau and Liebler (1992). Indeed this sort of analysis forms the basis for automated personality
115 analysis systems. Aran and Gatica-Perez (2013) focused on personality perception in a small group meeting
116 scenario. They extracted a set of multimodal features including speaking turn, pitch, energy, head and body
117 activity and social attention features. Thin slice analysis yielded the highest accuracy for *extroversion*,
118 while *openness* was better modelled by longer time scales. With regard to the related work in personality
119 computing, the closest approach was presented in Batrinca et al. (2016). In order to analyse the Big Five
120 personality traits, Batrinca *et al.* conducted a study where a set of participants were asked to interact with a
121 computer, which was controlled by an experimenter, and then a different set of participants were asked to
122 interact with the experimenter face-to-face to collaborate on completing a map task. In order to elicit the
123 participants personality traits, the experimenter exhibited four different levels of collaborative behaviors
124 from fully collaborative to fully non-collaborative. Self-reported personality traits were used to study the
125 manifestation of traits from audio-visual cues. In the human-machine interaction setting, their results
126 showed that 1) extroversion and neuroticism can be predicted with a high level of accuracy, regardless of
127 the collaboration modality; 2) prediction of the agreeableness and conscientiousness traits depends on the
128 collaboration modality; 3) openness was the only trait that cannot be modelled. In contrast to their findings
129 in the human-machine interaction setting, they showed that openness was the trait that can be predicted
130 with highest accuracy in the human-human interaction setting.

131 Applying such personality perception analysis to robot tele-operators has so far been limited. Perception
132 of tele-operator's personality is important not only in social interactions, but is also crucial where tele-
133 operated robots are used in a service capacity such as for elderly care Yamazaki et al. (2012), and search
134 and rescue Martins and Ventura (2009). In these settings, perception of the operator will effect system
135 utility for carrying out the desired service and achieving the desired outcome. In Celiktutan et al. (2016),
136 we showed that many of the aforementioned personality cues can be transmitted by a tele-presence robot.
137 We trained Support Vector Machine classifiers with a set of features extracted from participants' voice
138 and body movements. We found that the use of a robot avatar helps to discriminate between different
139 personality types (e.g., extroverted vs. introverted) better than audio-only mediated communication for
140 extroversion (65%) and conscientiousness (60%).

141 Studies with Mobile Remote Presence devices (MRPs) have briefly mentioned perceiving the operator's
142 personality Lee and Takayama (2011), but it has not been deliberately studied as we do here. There
143 are two studies that look directly at personality perception of tele-operators. Kuwamura et al. (2012)
144 examined an effect that they term *personality distortion*, demonstrated by reduction in internal consistency
145 of the personality questionnaire they used, for two different robot platforms and communication using
146 video. They use 3 tasks: (1) an experimenter talks freely with the participant, (2) a different experimenter
147 introduces and talks about themselves, and (3) a third experimenter interviews the participant. They only
148 observed *personality distortion* for one of the robot platforms, for *extroversion* in the interview task, and
149 for *agreeableness* in the introduction task. Using a single fixed person for each task, particularly members
150 of the experimental team who are aware of the goals of the study, greatly reduces the ecological validity of
151 their results. In contrast, here we use a large number of naïve targets performing naturalistic communication,
152 and conduct far more in-depth data analysis.

153 In a study with a tele-operated, highly humanlike robot, Straub et al. (2010) examined both how participant
154 tele-operators incorporate the fact that they are operating a robot into their presented identity, and how

155 interlocutors at the robot's location blend operator and robot identities. They used language analysis to
156 make their assessments. They observed that many operators pretended they themselves were a robot,
157 and interlocutors often referred to the operator as a robot. These behaviours are different from what
158 we typically observe with our tele-operation system, where most operators appeared to act naturally as
159 themselves Bremner et al. (2016b).

3 MATERIALS AND METHODS

160 We designed a two-stage experimental method for assessing changes in perceived personality that we used
161 in two studies. Firstly, a set of participants (targets) were recorded performing three communication tasks
162 in two conditions, directly visible on video camera (audio-visual condition) and communicating using the
163 tele-operated robot (tele-operated robot condition, also recorded on camera). This ensures we have a large
164 set of natural communication behaviours, and hence personality cues, for a range of personality types, that
165 can be viewed directly or when mediated by a robot.

166 In the second stage of the study, the recorded data was used to create a set of video clips for each target
167 in each communication condition. The video clips were pseudo-randomly assigned to a set of surveys in
168 such a way as to have one of each task and communication condition combinations present, with a given
169 target only appearing once in a given survey (i.e., communication condition was varied between surveys).
170 Each survey was viewed by a set of 10 judges, who after watching each clip assessed the personality of that
171 target. We used an online crowd-sourcing service to have the clips assessed. Employing judges via online
172 crowd-sourcing services has recently gained popularity due to its efficiency and practicality as it enables
173 collecting responses from a large group of people within a short period of time Biel and Gatica-Perez
174 (2013); Salam et al. (2016).

175 Personality was assessed by a questionnaire that aims to gather an assessment along the widely known Big
176 Five personality traits Vinciarelli and Mohammadi (2014). These five personality traits are *extroversion* (EX
177 - assertive, outgoing, energetic, friendly, socially active), *agreeableness* (AG - cooperative, compliant,
178 trustworthy), *conscientiousness* (CO - self-disciplined, organized, reliable, consistent), *neuroticism* (NE -
179 having tendency to negative emotions such as anxiety, depression or anger) and *openness* (OP - having
180 tendency to changing experience, adventure, new ideas). Each trait is measured using a set of items (the BFI-
181 10 Rammstedt and John (2007) with 2 per trait in the Solo Tasks Study, and the IPIP-BFM-20 Topolewska
182 et al. (2014) with 4 per trait in the Dyadic Tasks Study) scored on 10-point Likert scales. As well as being
183 assessed by external observers, each target completed the personality questionnaire for self assessment.

184 3.1 Tele-Operation System

185 In order to reproduce the gestures of targets on the NAO humanoid robot platform from Softbank
186 Robotics Gouaillier et al. (2009), we used a motion capture based tele-operation system. Previously we
187 have demonstrated the system to be capable of producing comprehensible gestures Bremner and Leonards
188 (2015, 2016). The arm motion of the targets is recorded using a Microsoft Kinect and Polhemus Patriot¹,
189 and used to produce equivalent motion on the robot. Arm link end points at the wrist, elbow and shoulder
190 are tracked, and were used to calculate joint angles for the robot so that its upper and lower arm links
191 reproduce human arm link positions and motion. This method ensures that joint coordination, and hand
192 trajectories are as similar as possible between the human and the robot within the constraints of the NAO

¹ Product of <http://polhemus.com/>

193 robot platform. Figure 1 shows a gesture produced by one of the targets, and the equivalent gesture on the
194 NAO.

195 3.2 Solo Tasks Study

196 3.2.1 Tasks

197 In the first study the three tasks performed by participants involved them performing directly to the
198 camera, i.e., solo, and were based upon a subset of tasks used by Borkenau et al. (2004). Each of the
199 tasks was framed as an interaction with the experimenter who stood beside the video camera used in the
200 recordings, and provided non-verbal feedback and prompt questions to ensure as natural communicative
201 behaviours as possible. Targets were instructed to speak for as long as they felt able, with a maximum time
202 of 2 minutes for each task. The majority of the targets talked for 30-60 s on each task, with occasional
203 prompts for missing information. Prior to performing tasks, we asked the targets to introduce themselves
204 and give some information about themselves, e.g., where they work, what they do, their family, etc. This
205 stage was purely to help naturalise the target to the experimental setting. It was not used to produce clips
206 for judge rating.

207 **Task 1 (Hobby):** This task asked targets to describe one of their hobbies, providing as much detail as
208 possible. Suggested detail included what their hobby involves, why they like it, how long have they been
209 doing it for, etc. Example personality cues we anticipated from this task include what targets have as their
210 hobby, and what detail and the depth of detail they provide while describing their hobby.

211 **Task 2 (Story):** This task is based on Murray's thematic apperception test (TAT), where the target is
212 shown a picture and is asked to tell a dramatic story based on a picture Murray (1943). They are asked
213 what is happening in the picture², what are the characters thinking and feeling, what happens before the
214 events in the picture and what happens after. The picture is purposely designed to be ambiguous so that
215 the target has the scope to interpret the picture as they see fit, and has to be creative in their story telling. It
216 is a projective test, where the details given by the target, and how they relate the actions of the characters,
217 provide cues about their personality.

218 **Task 3 (Mime):** This task required the targets to mime preparing and cooking a meal of their choice.
219 This was different from the mime task used by Borkenau et al. Borkenau et al. (2004), where targets had
220 to mime alternative uses for a brick. Our pre-tests showed little variability between targets for that task.
221 Instead, the chosen task gave the desired variability, and the gestures were better suited to performance on
222 the NAO robot. Which meal was selected, and the complexity of the mime, are example personality cues
223 we anticipated from this task.

224 3.2.2 Participants

225 26 participants were recorded as targets (16 female, Mean Age=30.85, SD=7.58), and gave written
226 informed consent for their participation, they were reimbursed with a £5 gift voucher for their time.
227 Recordings for 20 of the targets were used to create the clips used for judgements (6 targets were omitted
228 due to recording problems). The study was approved by the ethics committee of the Faculty of Environment
229 and Technology of The University of the West of England.

230 Clip ratings were undertaken by 143 judges recruited through the CrowdFlower online crowd sourcing
231 platform³. Judges were compensated 50 cents for annotating a total of four clips.

² Image used was <https://www.flickr.com/photos/bassclarinetist/>, used under creative commons licence.

³ CrowdFlower, a data enrichment, data mining and crowdsourcing company, <http://www.crowdfLOWER.com/>

232 3.2.3 Recordings

233 All tasks were recorded by one RGB video camera and the motion capture system used for tele-operation.
234 The recorded motion capture data was then used to produce robot mediated versions of the targets'
235 performances on the NAO robot using the aforementioned tele-operation system, which were also recorded
236 on video.

237 In addition to the audio-visual and tele-operated robot conditions, an audio only condition was created
238 using the audio from hobby and story tasks. Hence, each target had a total of 8 clips split over 3
239 communication conditions: 3 clips for the audio-visual condition, 2 clips for the audio-only condition,
240 and 3 clips for the tele-operated robot condition. This resulted in a total of 158 clips (two clips became
241 corrupted).

242 To avoid confusion, prompt questions were edited out of the clips. Further, for the few tasks where
243 performance exceeded 60 s, clips were edited to be close to this length as pre-tests showed a decrease in
244 the reliability of judgements with overly long clips. Mean clip duration was 50 s (SD=20 s).

245 The clips were split up into surveys each containing four clips: one of each task and one of the audio-only
246 clips, each of a unique target. Communication condition was pseudo-randomised across the three tasks in
247 each survey, but always contained at least one of each communication condition.

248 3.3 Dyadic Tasks Study

249 3.3.1 The Extended Tele-operation System

250 The tele-operation system was extended to enable interactive multi-modal communication. The first
251 addition made was a stereo camera helmet on the NAO robot, the images from which are displayed in an
252 Oculus Rift head mounted display (HMD). Coupled with using the Rift's inertial measurement unit to drive
253 the robot's head, meant the operator could see from the robots point of view, and their gaze direction and
254 head motion could be observed on the robot. Secondly we used a voice over IP communication system
255 to allow full duplex audio communication. Finally due to feedback from participants in the Solo Tasks
256 Study, we did not use the Polhemus Patriot in the Dyadic Tasks Study to make behaviours more natural;
257 importantly, wrist rotation was only really needed for the mime task in the Solo Tasks Study, and is less
258 important for normal gesturing. Figure 2 shows the tele-operation system and the setup during performance
259 of dyadic tasks in the tele-operation (TO) condition.

260 3.3.2 Tasks

261 In the second study the three tasks performed by participants involved interacting with a confederate, i.e.,
262 dyadic. A confederate was used to ensure that each participant had the same interactive partner, giving us a
263 measure of control over the interactions, while still seeming natural to the participants. The three selected
264 tasks were based on the suggestions in Funder et al. (2000) of having an informative task, a competitive
265 task and a cooperative task. The intention of these task types is that they each engender personality cues in
266 different ways.

267 The three tasks were briefly explained to the participant and the confederate together, and more detailed
268 written instructions were provided to be used during the experimental session. This was done to ensure
269 that the experimenters could leave the room for the participant and confederate to converse alone. The
270 two communication conditions (audio-visual and tele-operated robot) were performed sequentially, in a
271 pseudo-randomised order, in the same room. The audio-visual condition was recorded face-to-face, i.e.,
272 with both participant and confederate seated across a table from one another. In the tele-operated robot

273 condition the participant moved to an adjoining room where the tele-operation controls were located, while
274 the confederate sat at a table across from the robot.

275 **Task 1 (Informative):** Participants watched a clip from a Sylvester and Tweety cartoon, which they then
276 had to describe to the confederate. This is a task commonly used to examine gesturing Alibali (2001), as
277 describing the action filled cartoon often engenders gestures, which may be useful personality cues that
278 can be produced by the robot. Another key reason for this task choice was that all participants have the
279 same things to talk about: in the previously used hobby task several participants struggled to find much to
280 say without significant prompting. Two different Sylvester and Tweety cartoons were used, one for each
281 communication condition; cartoon assignment was randomised between conditions. We expected there to
282 be an abundance of gestural cues, as well as cues related to the participants' verbal behaviour (such as how
283 detailed the description was).

284 **Task 2 (Competitive):** The participants and the confederate played a memory based word game adapted
285 from the traditional *Grandmothers Trunk* game. The first player says "My Grandmother went on holiday
286 and she..." and adds something she did, accompanied by a gesture, the other player then repeats what the
287 first said and their gesture, and adds something else she did. Play continues alternating between players
288 who repeat the whole list of things and perform the gestures, adding a new thing each time, until one player
289 forgets something and that player loses. How they approach the competitive nature of the task, and the
290 actions they select are personality cues we expected from this task.

291 **Task 3 (Co-operative):** The participants and the confederate co-operated to put a set of 5 items into
292 utility order for surviving in a given scenario. There were two scenarios each with its own set of items,
293 surviving a ship wreck, and surviving a crash landing on the moon. One scenario was presented per
294 communication condition, and was randomly assigned. How agreement is reached, and how the task is
295 approached are the main cues we expect from this task.

296 3.3.3 Participants

297 30 participants were recorded as targets (13 female, Mean Age=25.01, SD=4.2), and gave written
298 informed consent for their participation, they were reimbursed with a £5 gift voucher for their time.
299 Recordings for 25 of the targets were used to create the clips used for judgements (5 targets were omitted due
300 to recording problems). The study was approved by the ethics committee of the University of Cambridge.

301 Clip ratings were undertaken by 250 judges recruited through the Prolific Academic online crowd
302 sourcing platform⁴. Each judge rated 6 clips and was compensated £2 for their time.

303 3.3.4 Recordings

304 In all tasks both the confederate and the participant were recorded by separate RGB video cameras.
305 The confederate was only recorded to obscure the fact that she was a confederate. In the tele-operated
306 robot condition a video camera recorded the robot instead of the participant. In order to produce videos of
307 identical length for all targets and tasks, the video clips were further edited to select a 60 s segment from
308 the beginning of the Informative task and from the end of Competitive and Co-operative tasks. This is in
309 line with suggestions by Carney et al. (2007b) for using clips of this length of a task to maximize consistent
310 judgement conditions for each target. Thus, each target had a set of three 60s clips for each of the two
311 communication conditions. One survey consisted of a pseudo-randomised set of 6 clips, 1 example of each
312 task in each communication condition, with unique targets in each clip. Additionally a practice clip of the

⁴ Prolific Academic online crowd sourcing platform, <https://www.prolific.ac/>

313 confederate was added to the start of all surveys to use as a measure of judge reliability, it also served to
 314 demonstrate her voice such that it could be ignored when she spoke during the target clips.

315 In Table 1, we summarised both studies in terms of number of participants, tasks, communication
 316 conditions and communicated cues.

Table 1. Summary of the conducted studies. AO: Audio-Only; Audio-Visual; TO: Tele-Operation.

Study	Num. of Participants	Tasks	Communication Conditions	Communicated Cues
Solo	26	Hobby, Story, Mime	AO, AV, TO	wrist, elbow, shoulder motion, wrist orientation
Dyadic	30	Informative, Competitive, Co-operative	AV, TO	wrist, elbow, shoulder motion; head motion; gaze direction

4 RESULTS AND ANALYSIS

317 To address the research questions introduced in Section 1, we analysed the level of agreement and the
 318 extent of shifts with respect to different communication conditions (e.g., audio-visual/AV, Audio-Only/AO,
 319 Tele-Operation/TO) and different tasks for each personality trait. We evaluated personality judgements to
 320 measure intra-/inter-agreement, self-other agreement and personality shifts as below.

- 321 • *Intra-judge Agreement.* Intra-judge agreement (also known as internal consistency) evaluates the
 322 quality of personality judgements based on correlations between different questionnaire items that
 323 contribute to measuring the same personality trait by each judge. We measured intra-judge agreement
 324 in terms of standardised Cronbach's α : $\alpha = \frac{K\bar{r}}{(1+(K-1)\bar{r})}$ where K is the number of the items ($K = 2$
 325 in the Solo Tasks Study, and $K = 4$ in the Dyadic Tasks Study) and \bar{r} is the mean of pairwise
 326 correlations between values assigned. The resulting α coefficient ranges from 0 to 1; higher values are
 327 associated with higher internal consistency and values less than 0.5 are usually unacceptable McKeown
 328 et al. (2012).
- 329 • *Inter-judge Agreement.* Inter-judge agreement refers to the level of consensus among judges. We
 330 computed the inter-judge agreement in terms of Intra-Class Correlation (ICC) Shrout and Fleiss (1979).
 331 ICC assesses the reliability of the judges by comparing the variability of different ratings of the same
 332 target to the total variation across all ratings and all targets. We used ICC(1,k) as in our experiments
 333 each target subject was rated by a different set of k judges, randomly sampled from a larger population
 334 of judges. ICC(1,k) measures the degree of agreement for ratings that are averages of k independent
 335 ratings on the target subjects.
- 336 • *Self-other Agreement.* Self-other agreement measures the similarity between the personality
 337 judgements made by self and others. We computed self-other agreement in terms of Pearson correlation
 338 and tested the significance of correlations using Student's t distribution. Pearson correlation was
 339 computed between the target's self-reported responses and the mean of the others' scores per trait.
- 340 • *Personality Shifts.* Personality shift refers to the extent to which people shifted from one personality
 341 class to another, in judges' perception, between AV and TO conditions. In order to measure shifts we
 342 first classified each target into low or high (e.g., *introverted* or *extroverted*) for each trait according to
 343 if their average judge rating for each task was above or below the mean for all targets in AV. For each

344 trait, each target was grouped according to their classification in both conditions, creating 4 groups (i.e.,
345 AV: high and TO: high, AV: high and TO: low, etc.). We presented these results in terms of contingency
346 tables and tested the significance using McNemar's test with Edwards's correction L.Edwards (1948).

347 In the following subsections, we present these results for each study (solo and dyadic) separately.

348 4.1 Solo Tasks Study

349 4.1.1 Elimination of Low-quality Judges

350 Although crowd-sourcing techniques have many advantages, identifying annotators who assign labels
351 without looking at the content (low-quality judges or spammers) is necessary to get informative results. As
352 a first measure we eliminated judges who incorrectly answered a test question about the content of the clips.
353 After this elimination mean-judges-per-clip was 7.9 (SD=1.5), with minimum judges-per-clip being 5.

354 To assess whether there remained further low-quality judges we calculated within-judge consistency for
355 the AV clips using Cronbach's α , which measures whether the values assigned to the items that contribute
356 to the same trait are correlated. The average value across all tasks was lower than we expected (less than
357 0.5), indicating some judges answer randomly. With no low-quality judges, we would expect values for the
358 AV clips greater than 0.5, i.e., in line with values reported in the literature for the BFI-10 with video clips
359 assessed by online judges Credé et al. (2012). We therefore used a judge selection method to remove these
360 additional low-quality judges. We used a ranking-based method based on pairwise correlations instead
361 of standard methods for outlier detection. For each clip, we calculated an average correlation score for
362 each judge from pairwise correlations (using all 10 questions in the BFI-10) with the remaining judges.
363 Judges with low correlation scores are deemed to be spammers. The judges were then ranked in order of
364 correlation score and the k highest ranked selected.

365 To evaluate the efficacy of this ranking procedure we calculated within-judge consistency results for
366 the AV clips for different judge numbers ranging from $k = 10$ (without elimination) to $k = 3$. These
367 values averaged over all tasks are presented in Figure 3-a. We further validated this by computing ICC
368 with varying number of judges, Figure 3-c. Selecting 5 judges per clip (based on pairwise comparisons)
369 was found to be sufficient to increase reliability to acceptable levels for the AV clips (greater than 0.5)
370 for all traits except for *openness*. We use 5 judges as it allows us to exclude all judges who failed the test
371 question while having the same number of judges for all clips (5 judges is common in this type of study,
372 e.g., Borkenau and Liebler (1992)).

373 4.1.2 Within-judge Consistency

374 Within-judge consistency was measured in terms of Cronbach's α . For the selected 5 judges per clip, the
375 detailed results with respect to different communication conditions and tasks are presented in Table 2-a,
376 where α values that indicate sufficient reliability for the BFI-10 (greater than 0.5, in line with values reported
377 in the literature Credé et al. (2012)) are highlighted in bold. To compare α values between communication
378 conditions we follow the method suggested by Feldt et al. (1987): 95% confidence intervals are calculated
379 for each α value, and if the value from one condition falls outside the confidence intervals from a condition
380 it is being compared to, this suggests it is significantly less consistent. Comparing AO with AV for the
381 hobby task, values for all traits, except for *agreeableness*, fall outside the 95% confidence intervals of the
382 AV values. Comparing TO with AV for the mime task, values for all traits, except for *conscientiousness*,
383 fall outside the 95% confidence intervals of the AV values. This indicates AV is found to be more consistent

384 as compared to AO for the hobby task (except for *agreeableness*) and TO for the mime task (except
385 for *conscientiousness*). No other comparisons indicate significant differences.

386 4.1.3 Between-judge Consistency

387 We computed between-judge consistency in terms of Intra-Class Correlation, ICC(1,k) proposed by Shrout
388 and Fleiss (1979), where $k = 5$. Our judge selection method uses the k most correlated judges so might
389 bias the ICC results (see Section 4.1.1). To evaluate this we calculated ICC for $k = (10, \dots, 3)$ for the
390 AV condition. Figure 3-b shows that, for *extroversion*, *conscientiousness* and *neuroticism*, ICC does not
391 change meaningfully as the number of judges varies, while selecting the 5 most correlated judges slightly
392 biases the results for *agreeableness* and *openness*.

393 The detailed results for the selected 5 judges per clip are presented in Table 2-b. We obtained significant
394 correlations for most traits in the AV condition, with values in the same range ($0.40 < ICC(1, k) < 0.81$)
395 as reported in the literature for online judges using a 10-item test ($0.42 < ICC(1, k) < 0.76$) Biel and
396 Gatica-Perez (2013). Fewer significant correlations were observed in the other communication conditions,
397 particularly in the story task for AO and the mime task for TO. *Extroversion* was the only trait that
398 consistently maintained correlation across conditions.

399 4.1.4 Self-other Agreement

400 We examined the extent to which judges agree with the target's self-assessment. Pearson correlations
401 between the self-ratings and the judge's ratings of conditions and tasks are reported in Table 2-c for the
402 selected 5 judges per clip. We observed that the judge's ratings bear a significant relation to the target's
403 self-ratings for *extroversion* only ($r = 0.24 - 0.44$ and $p < 0.05$). However, we did not obtain any
404 significant correlations in the TO condition (all $r < 0.2$ and $p > 0.05$).

405 4.1.5 Personality Shifts

406 We examined the extent to which people shifted from one personality class to another, in judges'
407 perception, between AV and TO conditions, in the hobby and story tasks for the selected 5 judges per clip.
408 We did not examine shifts involving AO or Mime task as the ICC scores indicated that personality ratings
409 in this condition would be too unreliable. These results are presented in Table 3 as 2x2 contingency tables.
410 To aid analysis we have also illustrated each shift as a proportional change (%) both from high to low
411 (HIGH2LOW) and from low to high (LOW2HIGH) in Figure 4 (see the figure on the left hand side).

412 We found a significant shift from high to low for *neuroticism* (70%). Note that the
413 corrected McNemar's test is very conservative in estimating significance, particularly for small
414 sample sizes. Although not statistically significant, we observed large shifts from low to high
415 for *extroversion* (56%), *conscientiousness* (67%) and *openness* (57%).

416 4.2 Dyadic Tasks Study

417 As in the Solo Tasks Study we assessed whether there existed low quality judges (spammers) in the judge
418 pool used for the Dyadic Tasks Study. To do so we repeated the same method that we used for the Solo
419 Tasks Study, where we evaluated ICC values, and used judge rating techniques to selectively remove judges.
420 These results are presented in Figure 3-b and -d. As we observed ICC values for the AV condition in line
421 with expectation with all judges included, and cannot observe large changes in the Cronbach's α values
422 and the ICC values, by excluding judges, we concluded that the judges were reliable. Hence, we present
423 the results for the Dyadic Tasks Study without eliminating any judges.

Table 2. Analysis of personality judgements across 3 communication conditions and 3 tasks. (a) Within-judge consistency in terms of Cronbach's α (good reliability > 0.80 is highlighted in bold); (b) Between-judge consistency in terms of ICC(1,k) (at a significance level of $*p < 0.05$, $**p < 0.01$, $***p < 0.001$); (c) Self-other agreement in terms of Pearson Correlation (at a significance level of $*p < 0.05$, $**p < 0.01$ and $***p < 0.001$).

	Audio-Visual (AV)				Audio-Only (AO)			Tele-Operation (TO)			
	Hobby	Story	Mime	All	Hobby	Story	All	Hobby	Story	Mime	All
<i>(a) Within-judge</i>											
EX	0.64	0.56	0.63	0.62	0.57	-0.15	0.34	0.61	0.39	0.19	0.47
AG	0.54	0.41	0.60	0.52	0.61	0.33	0.52	0.40	0.56	0.37	0.44
CO	0.47	0.60	0.54	0.55	0.50	0.21	0.39	0.54	0.56	0.57	0.55
NE	0.76	0.76	0.78	0.78	0.75	0.42	0.63	0.66	0.54	0.30	0.50
OP	-0.6	0.05	0.22	-0.04	-0.14	0.12	0.05	0.17	-0.24	-0.14	-0.07
<i>(b) Between-judge</i>											
EX	0.84***	0.81***	0.74***	0.81***	0.72***	0.51*	0.70***	0.72***	0.63**	-0.12	0.66***
AG	0.46*	0.61**	0.40	0.55***	0.25	-0.15	0.32	0.21	0.54**	-0.95	0.39**
CO	0.78***	0.67***	0.71***	0.72***	0.37	-0.10	0.22	0.32	0.65***	-0.35	0.36*
NE	0.80***	0.71***	0.55**	0.75***	0.57**	0.12	0.55***	0.70***	0.36	-0.56	0.44**
OP	0.12	0.67***	0.40	0.52***	0.49	0.40	0.55***	0.34	0.17	0.04	0.36*
<i>(c) Self-other</i>											
EX	0.34***	0.32**	0.26*	0.30***	0.44***	0.01	0.24***	0.12	-0.02	0.04	0.05
AG	0.04	0.13	0.04	0.07	0.28**	-0.05	0.12	0.08	-0.01	0.10	0.06
CO	-0.17	0.09	0.16	0.03	0.13	-0.13	0.01	0.05	0.16	-0.16	0.01
NE	0.00	-0.07	0.05	-0.01	0.07	0.09	0.07	0.02	-0.08	0.04	0.00
OP	0.06	0.03	0.00	0.03	0.10	0.04	0.07	0.16	0.07	0.03	0.09

Table 3. Contingency tables for each trait (at a significance level of $*p < 0.05$)

EX	TO: high	TO: low	AG	TO: high	TO: low	CO	TO: high	TO: low
AV: high	16	6	AV: high	16	11	AV: high	13	9
AV: low	10	8	AV: low	5	8	AV: low	12	6

NE	TO: high	TO: low	OP	TO: high	TO: low
AV: high	6	14*	AV: high	13	6
AV: low	1*	19	AV: low	12	9

424 4.2.1 Within-judge Consistency

425 Within-judge consistency was measured in terms of Cronbach's α . The detailed results with respect
 426 to different communication conditions and tasks are presented in Table 4-a, where α values that
 427 indicate sufficient reliability for the IPIP-BFM-20 (greater than 0.75, in line with values reported in
 428 the literature Credé et al. (2012)) are highlighted in bold. Values are above or close to good reliability
 429 (> 0.7) for all traits except for *neuroticism*. Comparing values across communication conditions we
 430 observe little difference, hence judges were able to make consistent trait evaluations when the robot is used
 431 for communication.

432 4.2.2 Between-judge Consistency

433 We computed between-judge consistency in terms of Intra-Class Correlation, ICC(1,k), where $k =$
434 10 Shrout and Fleiss (1979). The detailed results for the 10 judges per clip are presented in Table 4-
435 b. *Extroversion* and *openness* are the only traits with significant agreement across most tasks and both
436 conditions ($0.47 \leq ICC(1, k) \leq 0.85$ at a significance level of $p < 0.01$). Other traits vary between tasks
437 and conditions as to where significant agreement is achieved. A clearer picture can be gained from the all
438 task results, where it can be seen that agreement on *conscientiousness* deteriorates in the TO condition
439 relative to AV (a drastic drop from 0.61 to -0.26 over all tasks).

440 4.2.3 Self-other Agreement

441 We examined the extent to which judges agree with the target's self-assessment. Pearson correlations
442 between the self-ratings and the judge's ratings of conditions and tasks are reported in Table 4-c. Significant
443 agreement was found for *agreeableness* and *openness* across most tasks and both conditions ($r_{ag} = 0.75$
444 and $r_{op} = 0.71$ over all tasks), although agreement is much lower in the TO condition ($r_{ag} = 0.63$ and
445 $r_{op} = 0.46$ over all tasks). For *extroversion* and *neuroticism* agreement is much lower than for other traits,
446 and this is fairly consistent across conditions. Again we observe the larger difference across conditions
447 for *conscientiousness* ($r_{co} = 0.17$), with almost no significant agreement in the TO condition compared to
448 significant agreement across all tasks in the AV condition ($r_{co} = 0.31$).

449 4.2.4 Personality Shifts

450 We examined the extent to which people shifted from one personality trait classification to another, in
451 judges' perception, between AV and TO conditions for each task. These results are presented in Table 3 as
452 2x2 contingency tables. To aid analysis we have also illustrated each shift as a proportional change (%) both
453 from high to low (HIGH2LOW) and from low to high (LOW2HIGH) in Figure 4 (see the figure on the right
454 hand side). We found a significant shift from high to low for *agreeableness* (65%), *conscientiousness* (67%)
455 and *openness* (56%). Although not statistically significant, we observed a large shift from high to low
456 for *neuroticism* (57%).

5 DISCUSSION

457 In this section, we discuss our results, including comparisons with related work introduced in Section 2.
458 We present in depth discussion of meta-data (i.e., judge ratings, self ratings) in terms of intra/inter-
459 judge agreement, accuracy of judgements and personality shifts, with regard to different communication
460 conditions (i.e., AO: audio-only, AV: audio-visual, and TO: tele-operation) and different tasks (i.e., solo and
461 dyadic tasks). Note that in the majority of related works results were not directly comparable as personality
462 recognition accuracy is typically the reported metric, as opposed to agreement as used here; accuracy as
463 measured by comparing human responses with machine learning systems (e.g., Batrinca et al. (2016); Aran
464 and Gatica-Perez (2013)), or between self ratings and judge ratings (e.g., Funder (1995); Borkenau et al.
465 (2004)). Nevertheless, for which traits this reported accuracy is high or low helps provide some explanation
466 for our findings.

467 5.1 Intra-Judge Agreement

468 Consistency within judges for how each trait is judged (Table 2-a, Table 4-a) is used to address RQ1. In
469 both studies judges were sufficiently consistent in their trait ratings in the audio-visual condition (AV), with
470 the exception of *openness* in the Solo Tasks Study, and to a lesser extent *neuroticism* in the Dyadic Tasks

Table 4. Analysis of personality judgements across 2 communication conditions and 3 tasks. (a) Intra-judge consistency in terms of Cronbach's α (good reliability > 0.80 is highlighted in bold); (b) Inter-judge consistency in terms of ICC(1,k) (at a significance level of $*p < 0.05$, $**p < 0.01$, $***p < 0.001$); (c) Self-other agreement in terms of Pearson Correlation (at a significance level of $*p < 0.05$, $**p < 0.01$ and $***p < 0.001$).

	Audio-Visual (AV)				Tele-Operation (TO)			
	Informative	Competitive	Co-operative	All	Informative	Competitive	Co-operative	All
<i>(a) Within-judge</i>								
EX	0.85	0.87	0.85	0.87	0.84	0.85	0.84	0.86
AG	0.77	0.80	0.84	0.83	0.86	0.84	0.81	0.84
CO	0.71	0.75	0.77	0.74	0.76	0.70	0.72	0.73
NE	0.57	0.60	0.54	0.57	0.54	0.64	0.60	0.59
OP	0.78	0.82	0.87	0.85	0.75	0.79	0.85	0.81
<i>(b) Between-judge</i>								
EX	0.83***	0.84***	0.70***	0.85***	0.61***	0.78***	0.78***	0.82***
AG	0.18	0.21	0.58***	0.51**	0.08	0.35	0.37*	0.41*
CO	0.27	0.28	0.48**	0.61***	-0.24	-0.11	0.24	-0.26
NE	0.52**	0.53**	0.22	0.66***	0.38*	0.13	-0.35	0.46**
OP	0.21	0.67***	0.57***	0.51**	0.55**	0.47**	0.29	0.52**
<i>(c) Self-other</i>								
EX	0.29**	-0.12	-0.29**	-0.06	0.32**	0.21*	-0.15	0.18
AG	0.74***	0.73***	0.44***	0.75***	0.57***	0.65***	0.27**	0.63***
CO	0.22*	0.28**	0.31**	0.31**	-0.01	0.27**	0.14	0.17
NE	0.16	0.18	0.28**	0.24*	0.24*	0.19	0.07	0.23*
OP	0.68***	0.61***	0.17	0.71***	0.51***	0.37***	0.04	0.46***

Table 5. Contingency tables for each trait (at a significance level of $*p < 0.05$ and $***p < 0.001$)

EX	TO: high	TO: low	AG	TO: high	TO: low	CO	TO: high	TO: low
AV: high	31	5	AV: high	14	26***	AV: high	12	24*
AV: low	13	26	AV: low	5***	30	AV: low	10*	29

NE	TO: high	TO: low	OP	TO: high	TO: low
AV: high	16	21	AV: high	18	23*
AV: low	10	28	AV: low	10*	24

471 Study for us to conclude that the tasks and judges' behaviours were reliable. Batrinca et al. (2016) also
 472 reported a similar finding that openness was not modelled successfully in the human-machine interaction,
 473 whereas, in the human-human interaction setting, it was the only trait that could be predicted with a high
 474 accuracy over all collaboration tasks. In our case, the difference between the two studies with regards to
 475 consistent judgement of the *openness* trait indicates that cues for this trait may be more evident in dyadic
 476 tasks. Some researchers have suggested that one aspect of *openness* is intellect, where intellect incorporates
 477 the facets of intelligence, intellectual engagement and creativity DeYoung (2011), and the tasks in the
 478 Dyadic Tasks Study are more conducive to displaying these facets.

479 In the Solo Tasks Study there were some notable differences between the audio-only (AO) and the
 480 tele-operated robot (TO) conditions. For the hobby task, judges remained consistent in both the AO and
 481 TO conditions, indicating they were able to use audio cues to make judgements for this task, and robot

482 appearance had no effect on consistency. However, for the story task, judges were much less consistent
483 in the AO than in the AV condition, for all traits except for *agreeableness*. This is in contrast to the
484 tele-operated robot condition (TO), where they remained as consistent as in the AV condition. The only
485 additional cues available with the robot compared to audio only are gestures and appearance. The results
486 indicate that such cues are used to aid judgements in the same way that they do in the AV condition, though
487 their utility appears to be task dependent (only of apparent benefit in the story task). Importantly, the fact
488 that they are utilised provides good evidence that the robot is not simply ignored when making judgements.
489 Hence, the findings of high levels of agreement across both conditions in all tasks in the Dyadic Tasks
490 Study, indicate that in dyadic tasks the robot transmits sufficient cues to make judgements as consistently
491 as observing the target directly.

492 The use of gesture to aid personality judgements appears to be dependent on it accompanying speech, as
493 in the Solo Tasks Study ratings in the TO condition are far less consistent than in the AV condition for the
494 mime task. That is to say, gestures alone do not provide sufficient information for judging personality. This
495 was in contrast to what was reported by Aran and Gatica-Perez (2013), where the best results were achieved
496 when they used visual cues only for predicting personality traits, and using audio cues or combining them
497 with visual cues resulted in lower accuracy. This showed that either other behaviour cues not transmitted
498 by the robot are needed, or appearance cues are used which conflict with gesture cues in the TO condition.

499 Taking the results from both studies together it is apparent that judges are able to remain consistent in
500 their judgements of a given trait whether they are observing someone directly or their communication
501 relayed through a tele-operated robot. Indeed, where there are slight shifts in consistency between AV and
502 TO conditions they are not large; the one exception being for the mime task in the Solo Tasks Study. Hence,
503 each judge appears to formulate a relatively consistent evaluation of a given targets' personality traits based
504 on speech, gesture and appearance, combining them to assess each trait facet. This finding is in contrast
505 to Kuwamura et al. (2012) where they suggested small shifts in intra-judge consistency provided evidence
506 of robot appearance effects on personality perception. While in subsequent sections we do observe evidence
507 for effects of robot mediation on perception, we do not find such small shifts in intra-judge consistency
508 convincing in this regard.

509 5.2 Inter-Judge Agreement

510 Looking at inter-judge agreement results to address RQ2 (Table 2-b, Table 4-b), *extroversion* was the
511 only trait on which judges reached consensus in both studies, regardless of the communication condition,
512 and task (the mime task in the Solo Tasks Study being the one exception). This result is in line with the
513 widely accepted idea that *extroversion* is the easiest trait to infer upon others Barrick et al. (2000). Hence,
514 the strength of the available cues was sufficient to overcome any conflict between appearance, vocal, and
515 gesture based cues. Indeed it indicates that judges had a common set of interpretations for the available
516 cues.

517 On the other hand, where agreement was reached on *agreeableness*, *conscientiousness*,
518 and *neuroticism* for some tasks in the AV condition in each study, it had mostly deteriorated in
519 the TO condition, and the AO condition in the Solo Tasks Study. The clearest example of this is
520 for *conscientiousness* taking all three tasks together in the Dyadic Tasks Study (and to some extent in the
521 Solo Tasks Study as well), where agreement drastically deteriorated in the TO condition as compared to
522 the AV condition. As explained in Macrae et al. (1996), physical appearance based impressions (facial and
523 vocal features) are often used in the judgement of *conscientiousness*. In particular, low *conscientiousness* is
524 conveyed by a child-like face Macrae et al. (1996), which the face of the NAO robot can be considered to

525 have, and this may conflict with the vocal cues of the operator. *Neuroticism* is mainly related to emotions,
526 and *agreeableness* is related to trust, cooperation and sympathy Zillig et al. (2002), both of which it seems
527 reasonable to suggest judges might perceive as being low for a robot (particularly NAO with its lack of
528 facial expressions), again creating conflicts. It would appear that judges do not have a consistent manner
529 with which to resolve such conflicts.

530 Task based analyses in the Solo Tasks Study show that for *agreeableness* and *conscientiousness* the
531 story task provides sufficient cues for agreement to be maintained in the TO condition, whereas the hobby
532 task does so for *neuroticism*. As agreement being maintained in the TO condition indicates sufficient cues
533 to overcome appearance/behaviour conflicts, it is instructive to consider how those tasks might relate to
534 the traits. In telling the story, targets might demonstrate their morality, and relation to others, components
535 of *agreeableness* Zillig et al. (2002). How well structured and clear the story is could relate to facets of
536 the *conscientiousness* trait. The hobby task on the other hand might demonstrate how self-conscious a
537 person is about their hobby, a facet of *neuroticism* Zillig et al. (2002). While these two tasks might provide
538 some cues for facets of the traits for which consistency was not maintained, they appear to do so in a way
539 that conflicts with cues related to the robot.

540 We also compared, differences in agreement between the TO and AO conditions in the Solo Tasks
541 Study. Where there is agreement in TO for *agreeableness*, *conscientiousness* and *neuroticism*, we found it
542 was greatly reduced for *agreeableness* and *conscientiousness*, and to a lesser extent for *neuroticism*. This
543 provides further evidence that physical cues, be they behavioural or appearance based, are utilised in the
544 TO condition. Again, this appears to be dependent on the presence of speech: in the mime task for the Solo
545 Tasks Study judges were unable to provide a consistent rating for any trait in the TO condition, in contrast
546 to the consistent ratings for *extroversion*, *conscientiousness*, and *neuroticism* in the AV condition. A likely
547 reason for this observation is that without vocal cues there is an increased reliance on appearance based
548 cues, often based on stereotypes Kenny et al. (1994), and judges do not have consistent stereotypes relating
549 to robot appearance.

550 Batrinca et al. (2016) showed that the prediction of agreeableness and conscientiousness in the
551 human-machine interaction setting and the prediction of conscientiousness and neuroticism were
552 highly dependent on the collaboration task, where the extroversion trait was the only trait yielding
553 consistent results over all tasks in both settings. Similarly, our task based analyses in the Dyadic
554 Tasks Study show that in the AV condition, while the co-operative task provided a higher level
555 of agreement for *agreeableness* and *conscientiousness*, the competitive task yielded better results
556 for *neuroticism* and *openness*. Indeed, the results are somewhat expected given the nature of the tasks: the
557 co-operative task was to agree upon how to order five items in a survival scenario, in which participants
558 were expected to exhibit the *agreeableness* facet of personality; the competitive task was more related
559 to creativity and intelligence, that are strongly associated with *openness* Zillig et al. (2002). Though
560 agreement is lower, it is still maintained for *agreeableness* in the co-operative task and *openness* in the
561 competitive task in the TO condition. This indicates that in these cases, for at least some of the judges,
562 either the vocal cues override the visual cues, or movement cues are utilised (with the vocal cues).

563 Taken together, the findings from both studies indicate that the ability of judges to make judgements
564 based on a common interpretation of cues is affected not only by communication condition but is also
565 dependant on the task. While in some cases it is apparent that a particular task is conducive to providing
566 more verbal cues than another for a particular trait (as indicated by higher agreement, and inferred from the
567 literature), whether these override the physical cues in the TO condition is hard to predict. Indeed, whether
568 clear cues in the AV condition translate into agreement in the TO condition vary a great deal between all

569 tasks. Hence, it seems reasonable to suggest that whether inter-judge consistency is observed also depends
570 on how much appearance cues are utilised for a given task and trait, and thus how all the cues interact. This
571 complex interaction effect provides strong evidence that personality perception is likely to be altered when
572 communicating via a robot, and this depends on what cues are produced.

573 5.3 Accuracy of Judgements

574 In order to assess RQ3 we analysed the extent to which judge ratings correlated with self ratings provided
575 by target participants (Table 2-c and Table 4-c). In general in the Solo Tasks Study there was very little
576 correlation between self and other ratings. This is in contrast to previous findings where they found low,
577 but significant, self-other correlation (0.11 – 0.42) Carney et al. (2007a). The one exception to this was
578 self-other correlation for *extroversion* in the AV condition. This suggests that participant targets did not
579 present cues relating to their self-perception in the tasks we used, other than for *extroversion* which is
580 commonly reported as the trait with the most available cues. Audio cues were sufficient for this correlation
581 to be maintained in the hobby task in the AO condition, but not in the story task, or in either task in the TO
582 condition.

583 In contrast to the tasks used in the Solo Tasks Study, the tasks of the Dyadic Tasks Study resulted in
584 self-other agreement for *extroversion*, *agreeableness*, *conscientiousness*, and *openness* in the majority of
585 tasks for the AV condition. This indicates that the tasks we used in the Dyadic Tasks Study were better
586 at engendering more naturalistic behaviour, and hence personality cues than the tasks in the Solo Tasks
587 Study. Indeed, an important factor in thin slice personality analysis is how easy a person is to judge Funder
588 (1995), and people behaving more naturally produce better cues. However, despite these apparently better
589 cues, there was a large reduction in agreement for *conscientiousness*, *neuroticism*, and *openness* (and to
590 a lesser extent *agreeableness*) in the TO condition relative to the AV condition. This finding combined
591 with those of the Solo Tasks Study, suggests that there is a shift in the way personality cues are interpreted
592 caused by their interaction with the appearance of the robot, and the way non-verbal communication cues
593 are reproduced on it.

594 5.4 Personality Shifts

595 In order to address RQ4 we analysed the difference in perceived personality in terms of the occurrences
596 of personality shifts. We principally consider the results from the Dyadic Tasks Study as it provides the
597 more compelling evidence. The main reason for this assertion is that more naturalistic cues appeared
598 to be produced in the Dyadic Tasks Study (see previous section), and we consider such cues and their
599 interaction with the TO condition more ecologically valid. In addition, by being able to consider three
600 tasks rather than the two considered in the Solo Tasks Study we have increased statistical power. The shifts
601 we observed (Figure 4) provide evidence that cues related to the robots appearance are incorporated into,
602 or even override personality judgements based on speech. Indeed, this is somewhat to be expected given
603 that Behrend et al. (2012) observed that, in judgements of suitability, attractiveness of a graphical avatar
604 superseded qualities perceived in an interviewees words.

605 There are two likely causal factors in the perceived personalities being shifted, firstly human-based
606 physical appearance stereotypes (inferred from humanlike characteristics of the robot) might be applied,
607 secondly characteristics related to robots might be applied. Here we will discuss possible underlying causes
608 for the shifts observed in the Dyadic Tasks Study. In the case of *conscientiousness* and *neuroticism* a
609 childlike face, as the NAO might be considered to have, conveys low ratings for both traits Borkenau and
610 Liebler (1992); Macrae et al. (1996). Further, *conscientiousness* and *neuroticism* were also observed to

611 be influenced by face shape in graphical avatars Fong and Mar (2015), and as the NAO has a face
612 shape that differs from a human, hence this could lead to distortions in perceptions of these traits.
613 Additionally, *neuroticism* is mainly related to emotions Zillig et al. (2002), something which robots
614 are rarely considered to have. Also linked to emotions is *openness*, which combined with its other facets
615 of imagination and creativity, might also be reasonably expected to be low for a robot, which could also
616 be considered to have *hard facial linaments*, also linked to low *openness* Borkenau and Liebler (1992).
617 The NAO robot could also be considered male in appearance, and male avatars have been found to cue
618 for lower *conscientiousness* and *openness* Fong and Mar (2015). Low *agreeableness* is more difficult to
619 rationalise, but one facet is trustworthiness Zillig et al. (2002), and judges may have perceived using a
620 robot to communicate as less trustworthy. The vocal cues for *extroversion* appeared to be very strong, and
621 this might explain why little influence on this trait was observed.

622 An important thing to note from these findings is that people appear to be attributing personality
623 stereotypes to NAO for characteristics other than the *extroversion* trait which has been previously
624 examined Celiktutan and Gunes (2015); Aly and Tapus (2013); Park et al. (2012). Hence, in future
625 work in which a desired personality is to be expressed by an autonomous robot, its appearance based cues
626 must be considered alongside any behavioural cues expressed. We suggest that strong behavioural cues
627 may be required to overcome such stereotypes.

628 5.5 Conclusion

629 In this paper we have shown that judges are able to make personality trait judgements that are as consistent
630 with a robot avatar as when the same people are viewed on video in contrast to past work Kuwamura
631 et al. (2012). One possible reason for this difference in findings is that our tele-operation system allows
632 reproduction of some non-verbal communication cues on the robot which might improve the ease with
633 which judges can assess personality. Hence, we suggest that it is important for tele-presence systems to be
634 able to transmit non-verbal communication cues, whether this be actuation of physical systems, or large
635 enough screens on remote presence devices.

636 We have shown that the appearance of a tele-operated robot avatar influences how the personality of its
637 controller is perceived, i.e., robot appearance based personality cues are utilised along with cues in the
638 speech of the operators. Hence, the perceived personality of a tele-operator is shifted towards that related
639 to the robot's appearance. In light of these findings we suggest that robot avatar appearance and behaviour
640 be carefully considered relative to the person who will be controlling it, and this needs to be done on
641 an individual basis. Training of operators to produce clear cues, or having some cues appropriate to the
642 operator's personality autonomously generated, might allow some control of appearance effects.

643 Having the correct robot personality has been found to have a positive effect on interactions with
644 people Celiktutan and Gunes (2015); Aly and Tapus (2013); Park et al. (2012), and our findings also
645 have implications for such autonomous robot personality expression. It is important to consider what
646 appearance cues for personality a robot has, as we have observed humanlike personality inferences, and
647 whether the planned behavioural cues might conflict with them. Cues that work on one platform may not be
648 transferable to another. Additionally we suggest that future experiments on robots expressing personality
649 need to carefully consider tasks undertaken, as we observed that intra-judge agreement on personality
650 perception was highly task dependent.

651 5.6 Limitations and Future Work

652 While this paper provides evidence for how personality perception is affected for people tele-operating a
653 humanoid robot avatar, it has a number of limitations we hope to address in future work.

654 One area of limitation in our work relates to the movement capabilities of the NAO robot, and the inherent
655 differences with human movement capabilities. Although our previous work showed reproduced gestures
656 are comprehensible Bremner and Leonards (2015, 2016), there are clearly appreciable differences in the
657 way some movements are reproduced. Indeed, while these differences have limited affect on perceived
658 meaning, they likely contribute to the observed distortions in personality. The main limitations in this
659 regard are in elbow flexion, movement speed, and wrist and hand motion: the NAO elbow can only bend to
660 $\sim 90^\circ$, the main effect of which being a reduction in vertical travel of the hand for some gestures; humans
661 are capable of extremely rapid motions that the robot cannot match, consequently it will catch up as best it
662 can, but the usual response will be to not express some motions due to the method of motion processing;
663 wrist flexion and hand shape are clearly of utility in many gestures, and their absence (as well as wrist
664 rotation in study 2) restricts the expression of components of some gestures. These movement restrictions
665 are added to by limitations in the Kinect sensor and software processing: movements that result in hand
666 occlusions can lead to imprecision, as well as noise in the sensor data can lead to some added jitter on the
667 robot (though this is filtered as much as possible).

668 It is also important to note that robot operators had little to no awareness of the limitations of the robot as
669 none of them had prior experience with NAO, and when in control of it they could not observe its motion.
670 The only instruction given pertaining to system capabilities was to not to rest with the arms flat against the
671 body or behind the back as tracking would be lost. While this resulted in some initial poses that were a bit
672 unnatural (video of which was not used in the studies), participants soon reverted to 'normal' behaviour.
673 Indeed, qualitative comparison of participants in the dyadic study in each condition (video of participants
674 recorded while they were operating the robot allowed this) reveals little difference in gesturing behaviour
675 for the majority of participants. Exceptions were the two participants with prior experience working with
676 robots who moved more than they did face to face. In further work we aim to more closely examine the
677 data for any differences (which may be subtle), and if present test how they contribute to the observed
678 personality distortion effects.

679 In Celiktutan et al. (2016), our AV condition results showed that face gestures and head activity play
680 an important role in the recognition of the extroversion, agreeableness and conscientiousness traits. This
681 implies another limitation of the robotic platform used in this study. To convey the teleoperators personality
682 traits more accurately, the robot should portray head pose or facial activity together with audio and arm
683 gestures.

684 A further limitation is that there are some differences between our two studies, the Dyadic Tasks Study has
685 a slightly different design due to correcting issues we encountered in the Solo Tasks Study, making the
686 study comparison slightly less fair. In particular we addressed the issue with low quality judges, by utilising
687 a different recruiting platform which allowed us to recruit better quality judges, and thus did not require a
688 judge removal process. In the Solo Tasks Study the issues with low quality judges meant we used a judge
689 selection method based on the gathered responses. The procedure we used had a slight biasing effect on
690 the between-judge consistency (ICC) result for *agreeableness* and *openness*. This bias means that where
691 ICC values are not significant it is strong evidence that there is either a lack of cues or conflicting cues, as
692 even amongst the most agreeing judges consensus of opinion was not possible. Where there is significant
693 agreement, it indicates there are cues for that trait in the particular task and condition and some judges are

694 able to pick up on these cues. Indeed, Funder points out that there exists good and bad judges of personality
695 Funder (1995), and we suggest our selection method allowed us to bias toward good judges. This limits the
696 generalisability of our results to judges more adept at picking up on personality cues. By changing crowd
697 sourcing platforms we were able to remove the need for this selection process in the Dyadic Tasks Study.

698 In addition to recruiting better quality judges, we also utilised a larger personality questionnaire, making
699 our results more accurate, especially with regards to measuring intra-judge and inter-judge consistency.

700 In the work reported here it is not clear how different cues are utilised in the aforementioned personality
701 perception. Given that there was such high variability in affects of robot appearance dependent on the
702 task, it seems likely this is due to differences in use of audio and visual cues. Hence, we intend to analyse
703 in-depth the behaviours of targets relative to their judged personality for different tasks. To facilitate this
704 we aim to extend our work on automatic personality classification, which can extract and identify useful
705 cues automatically Celiktutan et al. (2016), and apply it to the recordings from the Dyadic Tasks Study. A
706 comparative cue analysis could not only allow us to gain a better understanding of the causes of personality
707 shifts, but could also be useful in synthesising robot personality behavioural cues.

AUTHOR CONTRIBUTIONS

708 PB: Substantial contributions to the conception and design of the work, the acquisition, analysis, and
709 interpretation of data. Drafting the work. Final approval of the version to be published. Agreement to
710 be accountable. OC: Substantial contributions to the conception and design of the work, the acquisition,
711 analysis, and interpretation of data. Drafting the work. Final approval of the version to be published.
712 Agreement to be accountable. HG: Substantial contributions to the design of the work, analysis, and
713 interpretation of data. Revising the work critically for important intellectual content. Final approval of the
714 version to be published. Agreement to be accountable.

ACKNOWLEDGMENTS

715 This work was funded by the EPSRC under its IDEAS Factory Sandpits call on Digital Personhood (Grant
716 Ref: EP/L00416X/1).

SUPPLEMENTAL DATA

717 Supplementary Material should be uploaded separately on submission, if there are Supplementary Figures,
718 please include the caption in the same file as the figure. LaTeX Supplementary Material templates can be
719 found in the Frontiers LaTeX folder

REFERENCES

- 720 Adalgeirsson, S. O. and Breazeal, C. (2010). MeBot: A robotic platform for socially embodied telepresence.
721 In *Proc. of Int. Conf. Human Robot Interaction*, pages 15–22. ACM/IEEE.
- 722 Alibali, M. (2001). Effects of Visibility between Speaker and Listener on Gesture Production: Some
723 Gestures Are Meant to Be Seen,. *Journal of Memory and Language*, 44(2):169–188.
- 724 Aly, A. and Tapus, A. (2013). A Model for Synthesizing a Combined Verbal and Nonverbal Behavior Based
725 on Personality Traits in Human-robot Interaction. In *Proc. of ACM/IEEE Int. Conf. on Human-Robot*
726 *Interaction*.

- 727 Aran, O. and Gatica-Perez, D. (2013). One of a Kind: Inferring Personality Impressions in Meetings. In
728 *Proc. of ACM Int. Conf. on Multimodal Interaction*.
- 729 Barrick, M. R., Patton, G. K., and Haugland, S. N. (2000). Accuracy of interviewer judgments of job
730 applicant personality traits. *Personnel Psychology*, 53(4):925–951.
- 731 Batrinca, L., Mana, N., Lepri, B., Sebe, N., and Pianesi, F. (2016). Multimodal personality recognition in
732 collaborative goal-oriented tasks. *IEEE Transactions on Multimedia*, 18(4):659–673.
- 733 Behrend, T., Toaddy, S., Thompson, L. F., and Sharek, D. J. (2012). The effects of avatar appearance on
734 interviewer ratings in virtual employment interviews. *Computers in Human Behavior*, 28(6):2128–2133.
- 735 Bevan, C. and Stanton Fraser, D. (2015). Shaking Hands and Cooperation in Tele-present Human-Robot
736 Negotiation. In *Proc. of Int. Conf. Human Robot Interaction*, pages 247–254. ACM/IEEE.
- 737 Biel, J. and Gatica-Perez, D. (2013). The YouTube Lens: Crowdsourced Personality Impressions and
738 Audiovisual Analysis of Vlogs. *Multimedia, IEEE Transactions on*, 15(1):41–55.
- 739 Borkenau, P. and Liebler, A. (1992). Trait inferences: Sources of validity at zero acquaintance. *J. of*
740 *Personality and Social Psychology*, 62(4):645–657.
- 741 Borkenau, P., Mauer, N., Riemann, R., Spinath, F. M., and Angleitner, A. (2004). Thin Slices of Behavior
742 as Cues of Personality and Intelligence. *J. of Personality and Social Psychology*, 86(4):599–614.
- 743 Bremner, P., Celiktutan, O., and Gunes, H. (2016a). Personality perception of robot avatar tele-operators.
744 In *The Eleventh ACM/IEEE International Conference on Human Robot Interaction, HRI '16*, pages
745 141–148.
- 746 Bremner, P., Koschate, M., and Levine, M. (2016b). Humanoid robot avatars: An 'in the wild' usability
747 study. In *RO-MAN*. IEEE.
- 748 Bremner, P. and Leonards, U. (2015). Efficiency of speech and iconic gesture integration for robotic and
749 human communicators - a direct comparison. In *Proc. of IEEE Int. Conf. on Robotics and Automation*,
750 pages 1999–2006. IEEE.
- 751 Bremner, P. and Leonards, U. (2016). Iconic gestures for robot avatars, recognition and integration with
752 speech. *Frontiers in Psychology*, 7:183.
- 753 Carney, D. R., Colvin, C. R., and Hall, J. A. (2007a). A thin slice perspective on the accuracy of first
754 impressions. *Journal of Research in Personality*, 41(5):1054–1072.
- 755 Carney, D. R., Colvin, C. R., and Hall, J. A. (2007b). A thin slice perspective on the accuracy of first
756 impressions. *Journal of Research in Personality*, 41(5):1054 – 1072.
- 757 Celiktutan, O., Bremner, P., and Gunes, H. (2016). Personality classification from robot-mediated
758 communication cues. In *25th IEEE International Symposium on Robot and Human Interactive*
759 *Communication (RO-MAN)*.
- 760 Celiktutan, O. and Gunes, H. (2015). Computational analysis of human-robot interactions through first-
761 person vision: Personality and interaction experience. In *24th IEEE International Symposium on Robot*
762 *and Human Interactive Communication (RO-MAN)*, pages 815–820.
- 763 Credé, M., Harms, P., Niehorster, S., and Gaye-Valentine, A. (2012). An evaluation of the consequences
764 of using short measures of the Big Five personality traits. *J. of personality and social psychology*,
765 102(4):874–88.
- 766 Daly-Jones, O., Monk, A., and Watts, L. (1998). Some advantages of video conferencing over high-quality
767 audio conferencing: fluency and awareness of attentional focus. *Int. Journal of Human-Computer*
768 *Studies*, 49(1):21–58.
- 769 DeYoung, C. D. (2011). Intelligence and personality. In Sternberg, R. J. and Kaufman, S. B., editors, *The*
770 *Cambridge handbook of intelligence*, pages 711–737. Cambridge University Press, New York.

- 771 Feldt, L. S., Woodruff, D. J., and Salih, F. A. (1987). Statistical Inference for Coefficient Alpha. *Applied*
772 *Psychological Measurement*, 11(1):93–103.
- 773 Fong, K. and Mar, R. A. (2015). What Does My Avatar Say About Me? Inferring Personality From Avatars.
774 *Personality and Social Psychology Bulletin*, 41(2):237–249.
- 775 Funder, D. C. (1995). On the accuracy of personality judgment: A realistic approach. *Psychological*
776 *review*, 102(4).
- 777 Funder, D. C., Furr, R. M., and Colvin, C. R. (2000). The riverside behavioral q-sort: A tool for the
778 description of social behavior. *Journal of personality*, 68(3):451–489.
- 779 Funder, D. C. and Sneed, C. D. (1993). Behavioral manifestations of personality: An ecological approach
780 to judgmental accuracy. *Journal of Personality and Social Psychology*, 64(3):479–490.
- 781 Gouaillier, D., Hugel, V., Blazevic, P., Kilner, C., Monceaux, J., Lafourcade, P., Marnier, B., Serre, J., and
782 Maisonnier, B. (2009). Mechatronic design of NAO humanoid. In *Proc of IEEE Int. Conf. on Robotics*
783 *and Automation*, pages 769–774. IEEE.
- 784 Hossen Mamode, H. Z., Bremner, P., Pipe, A. G., and Carse, B. (2013). Cooperative tabletop working
785 for humans and humanoid robots: Group interaction with an avatar. In *IEEE Int. Conf. on Robotics and*
786 *Automation*, pages 184–190. IEEE.
- 787 Kenny, D. A., Albright, L., Malloy, T. E., and Kashy, D. A. (1994). Consensus in interpersonal perception:
788 acquaintance and the big five. *Psychological bulletin*, 116(2):245–58.
- 789 Kristoffersson, A., Coradeschi, S., and Loutfi, A. (2013). A Review of Mobile Robotic Telepresence.
- 790 Kuwamura, K., Minato, T., Nishio, S., and Ishiguro, H. (2012). Personality distortion in communication
791 through teleoperated robots. In *Proc of IEEE Int. Symp. on Robot and Human Interactive*
792 *Communication*, pages 49–54. IEEE.
- 793 L.Edwards, A. (1948). Note on the correction for continuity in testing the significance of the difference
794 between correlated proportions. *Psychometrika*, 13(3):185–187.
- 795 Lee, M. K. and Takayama, L. (2011). Now, i have a body. In *Proc. of the conf. on Human factors in*
796 *computing systems*, page 33. ACM Press.
- 797 Macrae, C. N., Stangor, C., and Hewstone, M. (1996). *Stereotypes and Stereotyping*. The Guilford Press.
- 798 Martins, H. and Ventura, R. (2009). Immersive 3-d teleoperation of a search and rescue robot using a
799 head-mounted display. In *IEEE Conf. on Emerging Technologies Factory Automation (ETFA)*, pages
800 1–8.
- 801 McKeown, G., Valstar, M., Cowie, R., Pantic, M., and Schroder, M. (2012). The semaine database:
802 Annotated multimodal records of emotionally colored conversations between a person and a limited
803 agent. *IEEE Trans. Affect. Comput.*, 3(1):5–17.
- 804 Murray, H. A. (1943). *Thematic Apperception Test*. Harvard University Press.
- 805 Naumann, L. P., Vazire, S., Rentfrow, P. J., and Gosling, S. D. (2009). Personality judgments based on
806 physical appearance. *Personality & social psychology bulletin*, 35(12):1661–71.
- 807 O’Conaill, B., Whittaker, S., and Wilbur, S. (1993). Conversations Over Video Conferences: An Evaluation
808 of the Spoken Aspects of Video-Mediated Communication. *Human-Computer Interaction*, 8(4):389–
809 428.
- 810 Park, E., Jin, D., and del Pobil, A. P. (2012). The law of attraction in human-robot interaction. *International*
811 *Journal of Advanced Robotic Systems*, 9.
- 812 Rae, I., Takayama, L., and Mutlu, B. (2013). In-body experiences. In *Proceedings of the SIGCHI*
813 *Conference on Human Factors in Computing Systems - CHI '13*, pages 1921–1930, New York, New
814 York, USA. ACM Press.

- 815 Rammstedt, B. and John, O. P. (2007). Measuring personality in one minute or less: A 10-item short
816 version of the big five inventory in english and german. *J. of Res. in Personality*, 41(1):203 – 212.
- 817 Riggio, R. E. and Friedman, H. S. (1986). Impression formation: The role of expressive behavior. *Journal*
818 *of Personality and Social Psychology*, 50(2):421–427.
- 819 Salam, H., Celiktutan, O., Hupont, I., Gunes, H., and Chetouani, M. (2016). Fully automatic analysis of
820 engagement and its relationship to personality in human-robot interactions. *IEEE Access*, PP(99):1–1.
- 821 Shrout, P. and Fleiss, J. (1979). Intraclass correlations: Uses in assessing rater reliability. *Psychology Bull.*
- 822 Straub, I., Nishio, S., and Ishiguro, H. (2010). Incorporated identity in interaction with a teleoperated
823 android robot: A case study. In *Proc of Int. Symp. in Robot and Human Interactive Communication*,
824 pages 119–124. IEEE.
- 825 Tang, A., Boyle, M., and Greenberg, S. (2004). Display and presence disparity in Mixed Presence
826 Groupware. In *Proc. of Australasian User Interface Conf.*, pages 73–82. Australian Computer Society,
827 Inc.
- 828 Topolewska, E., Skiminia, E., Strus, W., CIECIUCH, J., and ROWINSKI, T. (2014). The short ipip-bfm-20
829 questionnaire for measuring the big five. *Annals Psychology*, 2(XVII):385–402.
- 830 Vinciarelli, A. and Mohammadi, G. (2014). A Survey of Personality Computing. *IEEE Trans. on Affective*
831 *Computing*.
- 832 Wang, Y., Geigel, J., and Herbert, A. (2013). Reading Personality: Avatar vs. Human Faces. In *Proc. of*
833 *HAC Conf. on Affective Computing and Intelligent Interaction*, pages 479–484. IEEE.
- 834 Yamazaki, R., Nishio, S., Ogawa, K., and Ishiguro, H. (2012). Teleoperated android as an embodied
835 communication medium: A case study with demented elderlies in a care facility. In *RO-MAN*, pages
836 1066–1071. IEEE.
- 837 Zillig, L. M. P., Hemenover, S. H., and Dienstbier, R. A. (2002). What do we assess when we assess a
838 big 5 trait? a content analysis of the affective, behavioral, and cognitive processes represented in big 5
839 personality inventories. *Personality and Social Psychology Bulletin*, 28(6):847–858.



Figure 1. Snapshots from the Solo Tasks Study. Left hand side: a target perceived to be *extroverted* by judges. Right hand side: a target perceived to be *introverted* by judges.

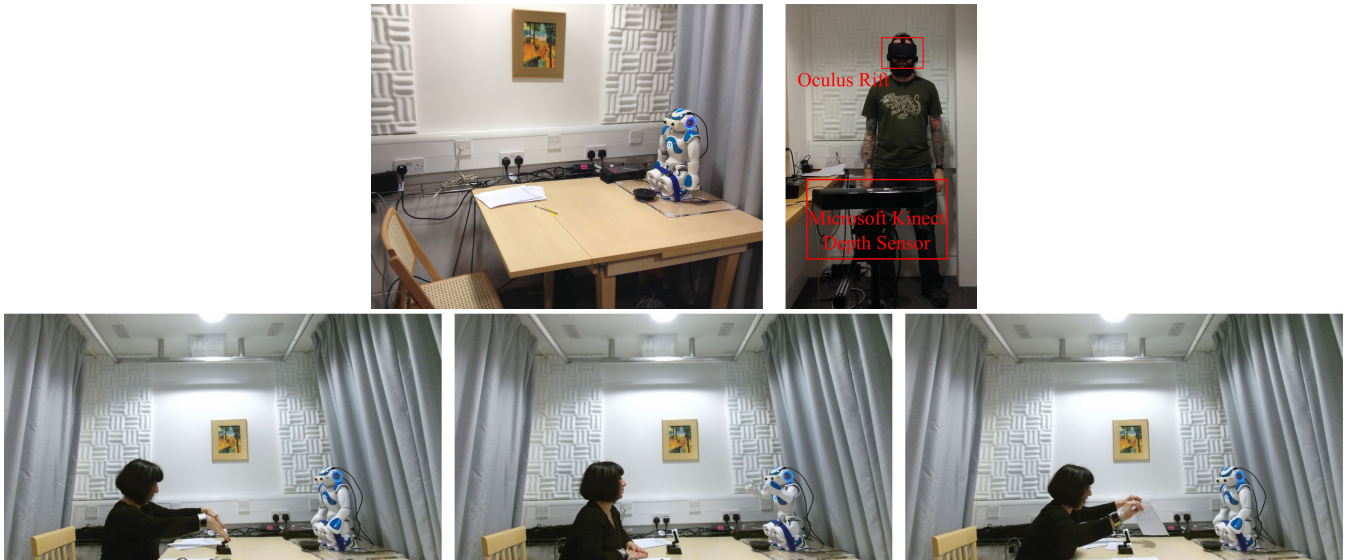


Figure 2. Snapshots from the Dyadic Tasks Study. Upper row: Illustration of tele-operation (TO) room and interaction room. Lower row: Snapshots from the dyadic interaction sequences.

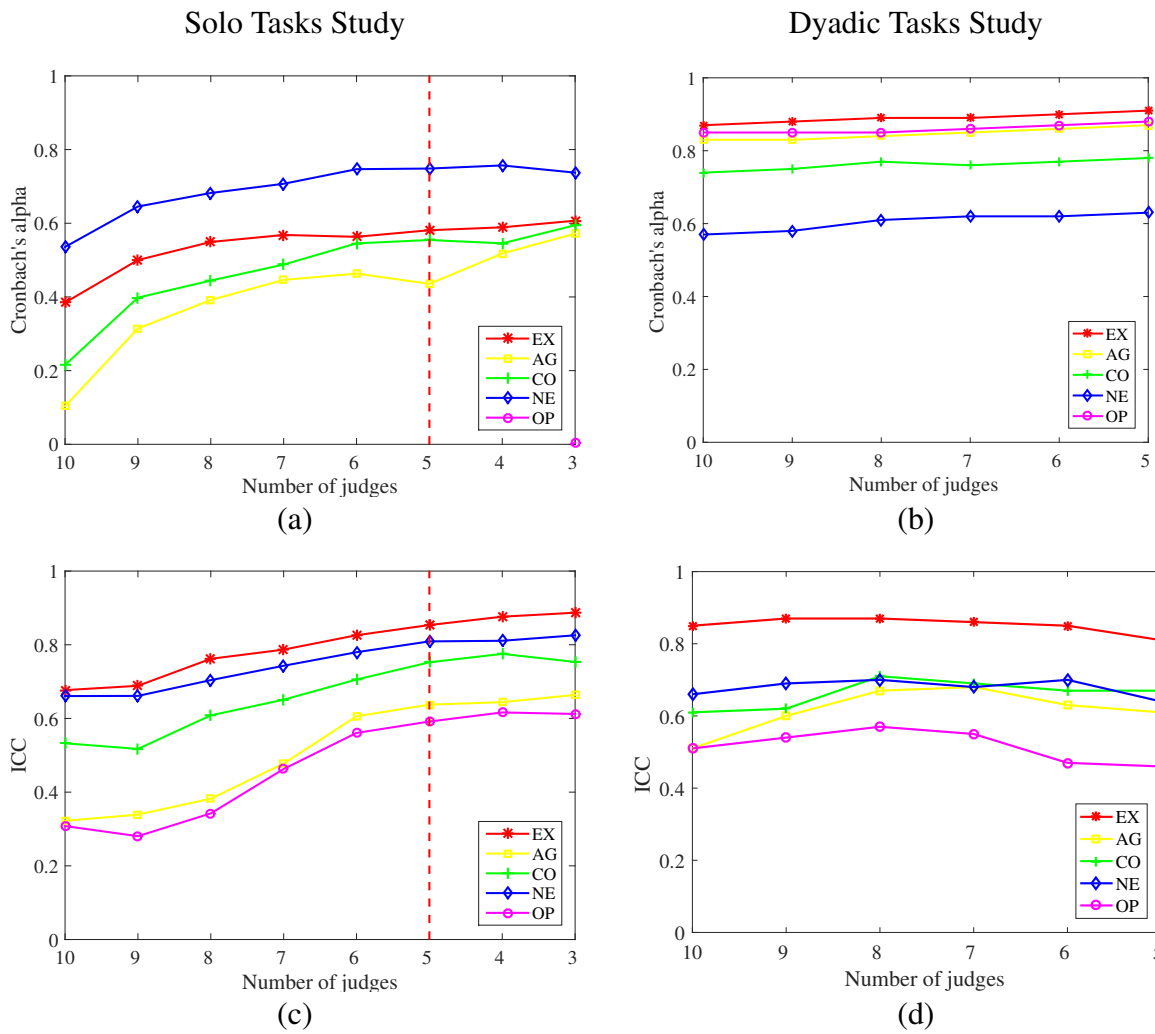


Figure 3. Changes in Cronbach's α values (a-b) and ICC values (c-d) as a function of number selected judges (k) for different traits in the AV communication condition for Solo Tasks Study (a-c) and Dyadic Tasks Study (b-d).

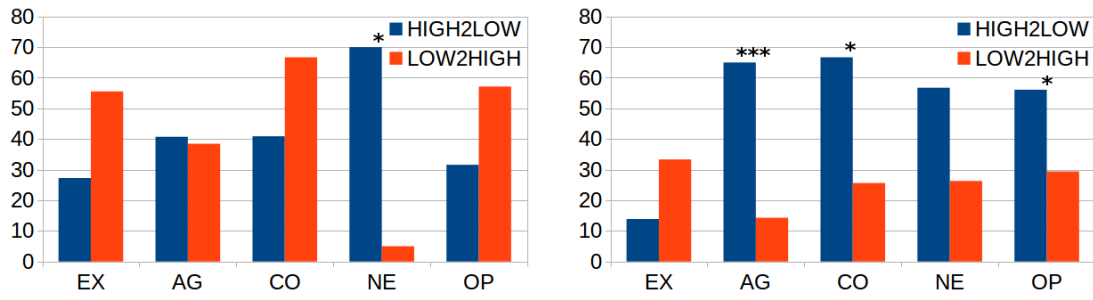


Figure 4. Amount of shifts (%) from high to low (HIGH2LOW) and from low to high (LOW2HIGH) (* : $p < 0.05$, *** : $p < 0.001$) between AV and TO: solo tasks (left hand side) versus dyadic tasks (right hand side).