
To Err is Robot: How Humans Assess and Act Towards an Erroneous Social Robot

Nicole Mirnig^{1*}, Gerald Stollnberger¹ Markus Miksch¹ Susanne Stadler¹

Manuel Giuliani² and Manfred Tscheligi^{1,3}

¹ *Center for Human-Computer Interaction, University of Salzburg, Austria*

² *Bristol Robotics Laboratory, University of the West of England, United Kingdom*

³ *Center for Technology Experience, Austrian Institute of Technology, Vienna, Austria*

Correspondence*:

Nicole Mirnig, Center for Human-Computer Interaction, University of Salzburg,
Jakob-Haringer-Straße 8/Techno 5, 5020 Salzburg, Austria
nicole.mirnig@sbg.ac.at

2 ABSTRACT

3 We conducted a user study for which we purposefully programmed faulty behavior into a robot's
4 routine. It was our aim to explore if participants rate the faulty robot different from an error-free
5 robot and which reactions people show in interaction with a faulty robot. The study was based on
6 our previous research on robot errors where we detected typical error situations and the resulting
7 social signals of our participants during social human-robot interaction. In contrast to our previous
8 work, where we studied video material in which robot errors occurred unintentionally, in the herein
9 reported user study, we purposefully elicited robot errors to further explore the human interaction
10 partners' social signals following a robot error. Our participants interacted with a human-like NAO,
11 and the robot either performed faulty or free from error. First, the robot asked the participants
12 a set of predefined questions and then it asked them to complete a couple of LEGO building
13 tasks. After the interaction, we asked the participants to rate the robot's anthropomorphism,
14 likability, and perceived intelligence. We also interviewed the participants on their opinion about
15 the interaction. Additionally, we video-coded the social signals participants showed during their
16 interaction with the robot as well as the answers they provided the robot with. Our results show
17 that participants liked the faulty robot significantly better than the robot that interacted flawlessly.
18 We did not find significant differences in people's rating of the robot's anthropomorphism and
19 perceived intelligence. The qualitative data confirmed the questionnaire results in showing that
20 although the participants recognized the robot's mistakes, they did not necessarily reject the
21 erroneous robot. The annotations of the video data further showed that gaze shifts (e.g., from and
22 object to the robot or vice versa) and laughter are typical reactions to unexpected robot behavior.
23 In contrast to existing research, we assess dimensions of user experience that have not been
24 considered so far and we analyze the reactions users express when a robot makes a mistake.
25 Our results show that decoding a human's social signals can help the robot understand that there
26 is an error and subsequently react accordingly.

27 **Keywords:** social human-robot interaction, robot errors, user experience, social signals, likeability, faulty robots, error situations,

28 **Pratfall Effect**

1 INTRODUCTION

29 Social robots are not yet in a technical state where they operate free from errors. Nevertheless, most research
30 approaches act on the assumption of robots performing faultlessly. This results in a confined standpoint,
31 in which the created scenarios are considered as gold standard. Alternatives resulting from unforeseeable
32 conditions that develop during an experiment are often not further regarded or simply excluded. It lies
33 within the nature of thorough scientific research to pursue a strict code of conduct. However, we suppose
34 that faulty instances of human-robot interaction (HRI) are nevertheless full with knowledge that can help us
35 further improve the interactional quality in new dimensions. We think that because most research focuses
36 on perfect interaction, many potentially crucial aspects are overlooked.

37 Research that is specifically directed at exploring erroneous instances of interaction could be useful to
38 further refine the quality of HRI. For example, a robot that understands that there is a problem in the
39 interaction by correctly interpreting the user's social signals, could let the user know that it understands the
40 problem and actively apply error recovery strategies. Knowing the severity of an error, could further be
41 helpful for the robot in finding the adequate corrective action.

42 Since robots in HRI are social actors, they elicit mental models and expectations known from human-
43 human interaction (HHI), Lohse (2011). One aspect we know from HHI is that imperfections make human
44 social actors more likeable and more believable. The psychological phenomenon *Pratfall Effect* states that
45 people's attractiveness increases when they commit a mistake. Aronson et al. (1966) suggest that superior
46 people may be viewed as superhuman and distant while a mistake would make them seem more human.
47 Similarly, one could argue that robots are often seen as impeccable, since this is how they are presented in
48 the media, Bruckenberger et al. (2013). Especially people who have not interacted with robots themselves
49 build their mental models and expectations about robots from those media. Moreover, experience with
50 technology in general is mostly based on interaction with consumer products, such as smartphones or TVs.
51 Those products are very common and need to work more or less error-free in order to get accepted on the
52 market. For example, a TV which has problems in sound will not survive long on the market. People expect
53 technology they paid for to work without errors. What makes the interaction with social robots different, is
54 that a TV is not seen as a social actor, in contrast to a social robot. This might result in people assuming
55 robots to be without fail which makes them likewise seem distant (*Pratfall Effect*). Robots that commit
56 errors, on the other hand, could then be viewed as more human-like and, in subsequence, more likeable.
57 With their study on an erroneous robot in a competitive game-play scenario Ragni et al. (2016) provided
58 additional evidence that people consider robots in general as competent, functional, and intelligent.

59 In our effort to embrace the imperfections of social robots and create more believable robot characters,
60 we propose to specifically explore faulty robot behavior and the social signals humans show when a robot
61 commits a mistake. The term social signal is used to describe verbal and non-verbal signals that humans
62 use in a conversation to communicate their intentions. Vinciarelli et al. (2009) argue that the ability to
63 recognise social signals is crucial to mastering social intelligence. It is our long-term goal to enable robots
64 to communicate about their errors and deploy recovery strategies. To achieve this ambitious goal, more
65 general knowledge about robot errors is required. We report on a user study where we purposefully elicited
66 faulty robot behavior.

67 Our user study is based on our previous research where we analyzed an extensive pool of video data
68 showing social HRI instances where the robot made an error. The videos covered a variety of scenarios
69 in different contexts, different robots, and a multitude of social signals. The robot errors happened
70 unintentionally and, thus, the data created a sound basis for studying the nature of error situations. We

71 found that there are two different kinds of robot errors, i.e., *social norm violations* and *technical failures*
72 Giuliani et al. (2015), for which human interaction partners respond with typical social signals, Mirnig
73 et al. (2015). A social norm violation means that the robot's actions deviate from the underlying social
74 script, that is the commonly known interaction steps a certain situation is expected to take. For example, a
75 participant orders a drink from a bartender robot, the robot signals it has understood but then asks again for
76 the participant's order. A technical failure means that the robot experiences a technical disruption that is
77 perceived as such by the user. For example, a robot picks up an object but then loses it while grasping. From
78 an expert perspective all robot errors might be considered as technical failures. Since, we are interested in
79 the human perception of robot errors, we distinguish error types from how a human most likely perceives
80 error events.

81 With the user study presented in this paper, we expand our previous research in purposefully eliciting
82 robot errors and researching the resulting social signals of the human interaction partners. We measured
83 how users perceive a robot that makes errors during interaction (social norm violations and technical
84 failures) as compared to a robot operating free from errors.

85 The directed exploration of robot errors in social interaction is a new and upcoming topic. The HRI
86 research community has reported first results on exploratory user studies. For example, Salem et al. (2015)
87 conducted an experiment with an erroneous robot. The researchers measured how the robot's behavior
88 influenced how the participants rated its trustworthiness and reliability. They also measured if robot errors
89 affect the task performance. The researchers found that while participants rated the correctly behaving
90 robot as significantly more trustworthy and reliable, the fact that a robot performs correctly or faulty did
91 not influence the objective task performance.

92 In an earlier work, Salem et al. (2013) researched the effect of speech and gesture congruence on perceived
93 anthropomorphism, likability, and task performance. In their experiment, a robot either spoke only, spoke
94 while making congruent co-verbal gestures, or spoke while making incongruent co-verbal gestures. The
95 researchers found that congruent co-verbal gesturing makes a robot appear more anthropomorphic, and
96 more likeable. This effect was even stronger for incongruent co-verbal gesturing. However, incongruent
97 co-verbal gesturing resulted in a lower task performance. Following our line of argumentation, such
98 incongruent behavior violates the human *social script*, as humans do not expect incongruent messages from
99 different modalities in everyday interactions. Therefore, incongruent multimodal robot behavior results in
100 a *social norm violation*. Ragni et al. (2016) report similar effects. The researchers performed a study in
101 which a human and a robot competed against each other in a reasoning task and a memory task. During
102 the interaction, the robot either performed with or without errors. While participants rated the faulty robot
103 as less competent, less reliable, less intelligent, and less superior than the error-free robot, participants
104 reported having enjoyed the interaction more when the robot made errors. However, the task performance
105 was significantly lower in the faulty robot condition.

106 Gompei and Umemuro (2015) investigated how a robot's speech errors influenced how familiar and
107 sincere it was rated. The researchers found that speech errors made early in an interaction might lower the
108 robot's sincerity rating. However, speech errors that are introduced later in the interaction might increase
109 the robot's familiarity. Short et al. (2010) investigated people's perception when playing rock-paper-scissors
110 with a robot that either played fair, cheated verbally by announcing a different hand gesture, or cheated
111 with its actions by changing the hand gesture. The researchers found that a cheating robot resulted in a
112 bigger social engagement, in comparison to one which plays fair. They stated, that the results suggest that
113 participants showed more verbal social signals to the robot that cheated. Participants were surprised by
114 the cheating behavior of the robot, although verbal cheating was perceived as malfunction, while cheating

115 through action was perceived as deliberate cheating behavior. These findings support our assumption, that
116 through unexpected behavior, people see a robot as a more social actor and that unexpected behavior might
117 be interpreted as erroneous behavior.

118 In an online survey, Lee et al. (2010) found that when a service robot made a mistake, this has a strong
119 negative impact on people's rating of the service quality and the robot itself. However, when the robot
120 deployed a recovery strategy, both the rating of the service and the rating of the robot improved. The
121 researchers deployed different recovery strategies and found that all of them increased the ratings of the
122 robot's politeness. A robot which apologized for its mistake was seen more competent, people liked it more
123 and felt closer to it, and a robot offering compensation for its mistake (such as a refund) was rated to be of
124 more satisfying service quality but participants were hesitant to use the robot again. Whereas, an apology
125 and a recovery strategy of offering options was perceived to foster re-use likelihood. In a related online
126 survey, Brooks et al. (2016) explored people's reactions to the failure of an autonomous robot. In the survey,
127 participants were asked to assess situations where an autonomous robot experienced different kinds of
128 failures that affected a human interacting with it. They found that people who saw an erroneous robot rated
129 it rather negatively on a series of items (i.e., How satisfying, pleasing, disappointing, reliably, dependable,
130 competent, responsible, trustworthy, risky to use is the robot?), while people who experienced a robot
131 without failure rated it positively. When the erroneous robot deployed mitigation strategies to overcome
132 the error either by prompting human intervention or by deploying a different approach, people's ratings
133 towards the erroneous robot became less negative. However, the amount the strategy influenced peoples
134 reaction depended on the kind of task, the severity of the failure, and the risk of the failure.

135 To enable a robot to generate help requests in case of an error situation, Knepper et al. (2015) developed
136 their inverse semantics algorithm. It allows the robot to phrase precise requests that specify the kind of help
137 that is needed. The researchers evaluated their algorithm in a user study and found that participants preferred
138 the precise request over high level, general phrasings. While in their approach errors are recognized through
139 the robot's internal state and the environment (e.g., the robot is supposed to pick up an object which it can
140 visually detect, but the object is out of its reach), we envision an approach where the robot can additionally
141 detect an error through its human interaction partner's social signals. For example, Gehle et al. (2015)
142 explored gaze patterns of human groups upon unexpected robot behavior in a museum guide scenario. They
143 found that groups of visitors responded to unexpected robot behavior with stepwise gaze coordination,
144 applying different modes of gaze constellation. Unexpected robot behavior is likely to conflict with the user
145 expectations about the adequate *social script* in a certain situation. Therefore, unexpected robot behavior
146 can lead to a social norm violation. A deviation from the *social script* resulted in a different strategy in the
147 human gaze coordination (social signals). Hayes et al. (2016) performed a user study in which participants
148 were instructed to teach a dance to a robot. They explored how humans implicitly responded when the
149 robot made a mistake. The authors used a very small sample in their explorative study and did not provide
150 a statistical analysis of their descriptive results.

151 Our approach extends the existing findings in several dimension. While the errors in Ragni et al. (2016)
152 were based on errors from HHI, the errors we used were modelled based on data from HRI. Our work and
153 Ragni et al. (2016) further cover different aspects: (a) their errors were task-related, ours non task-related;
154 (b) they covered the cognitive ability of the robot and we dealt with socially (in)appropriate robot behavior
155 and more general soft- and hardware problems; (c) they assessed the overall enjoyment of the interaction
156 and users' task performance, while we looked into the interconnectedness of likability, anthropomorphism,
157 and intelligence. We chose to examine these factors since they are commonly used and accepted measures
158 in the HRI domain. We were especially interested in likability as it contributes to the overall user experience

159 and it may foster technology acceptance. Since erroneous behavior potentially compromises intelligence
160 ratings, we were also interested in exploring if our robot's mistakes make them seem less intelligent. In the
161 light of the Pratfall Effect, we wanted to see if the robot's anthropomorphism level is influenced by the fact
162 that it makes or does not make mistakes.

163 The related literature shows that the importance of exploring robot errors has been recognized. We extend
164 the state of the art with our data-driven approach by systematically analyzing specific kinds of errors and
165 their effects on the interaction experience, as well as the users' reactions to those errors (i.e., social signals).

2 MATERIAL AND METHODS

166 We set up a Wizard of Oz (WOz) user study to specifically explore robot errors. A human and a robot
167 interacted with each other in two verbal sessions. The first session was a verbal interview where the robot
168 asked a few questions to the participant. The second session was a LEGO task, where the robot invited the
169 participant to build a few simple objects. We chose this setup in order to re-enact the verbal context of the
170 related work, Giuliani et al. (2015); Mirnig et al. (2015). In addition, the interview session enabled us to
171 collect qualitative data on the participants' opinions which we included in our data analysis.

172 The user study was performed between subjects, with each participant taking part in one of the following
173 two conditions: (a) *no error* (baseline - the robot performs error-free), and (b) *error* (experimental condition
174 - the robot commits eight errors over the entire interaction). To base the user study on the previous findings
175 from Giuliani et al. (2015) and Mirnig et al. (2015), we programmed the robot to commit two social norm
176 violations and two technical failures in each session. Based on our previous research, we defined these
177 two types of error as the typical mistakes robots make in HRI. Therefore, we suppose that an interaction
178 including these error types would be perceived as plausible. The complexity, severity, and risk-level of
179 the induced errors were chosen in alignment with our scenario. Naturally, different scenarios will entail
180 other errors, different severity and risk-levels. For example, Robinette et al. (2015) investigated faulty
181 behavior of robots in safety critical situations. They simulated erroneous behavior of an emergency guiding
182 robot that helps people to escape from a dangerous zone. They found that after the first error of the robot,
183 people's attitude toward the robot decreased significantly. However, the decision to follow the robot in a
184 follow-up interaction was not affected by their decreased attitude.

185 2.1 Hypotheses

186 As discussed in the previous sections, it is known that humans often base their expectations about robots
187 on how robots are portrayed in the media. Since the media present robots frequently as perfect entities, we
188 assume that social robots making errors negatively influence how their human interaction partners perceive
189 them. Based on the findings on faulty robot actions in HRI as discussed so far, we have postulated the
190 following hypotheses for our user study:

191 H1: A robot that *commits errors* during its interaction with humans, is perceived as *more likeable* than a
192 robot that performs flawlessly.

193 H2: A robot that *commits errors* during its interaction with humans, is perceived as *more anthropomorphic*
194 than a robot that performs flawlessly.

195 H3: A robot that *commits errors* during its interaction with humans, is perceived as *less intelligent* than a
196 robot that performs flawlessly.

197 2.2 User Study Design

198 For the WOz user study the participants were asked to interact with a NAO robot¹. We set the interaction
199 up in two sessions. During the first session, the robot asked a set of predefined questions to the participant
200 in order to restrict the thematic dimension of the conversation. During the second session, the robot invited
201 the participant to perform some tasks using LEGO bricks.

202 In the interview session the robot asked ten questions to the participant. The first three questions were
203 meant to make the participant familiar with the situation and to create a comfortable atmosphere. For this
204 reason, they were always presented in the same order and they never contained an error. The subsequent
205 seven questions were asked in random order and four out of seven questions contained errors in the *error*
206 condition.

207 In the LEGO session the participant had to (dis-)assemble LEGO bricks according to the robot's
208 instructions. The first two tasks were assigned in the same order for all participants and they did not contain
209 errors. The subsequent eight tasks were assigned in random order and four out of eight tasks contained
210 errors in the *error* condition.

211 The interview session lasted for an average of 3 minutes and 37 seconds ($SD = 59$ seconds), the LEGO
212 session about 8 minutes and 14 seconds ($SD = 1$ minute and 54 seconds). We decided for this two-part
213 setup to keep the participants entertained with a diversified scenario. The two-part setup provided us also
214 with the possibility to introduce a greater variety of errors and to achieve a higher number of errors in total.

215 The user study was performed in the User Experience and Interaction Experimentation Lab at the Center
216 for Human-Computer Interaction at the University of Salzburg. The robot was wizarded from a researcher
217 seated behind a bookshelf so that the wizarding was not obvious to the participant. A second researcher,
218 likewise seated behind the bookshelf, controlled the video recording. During the entire interaction the
219 participants stood adverse to the NAO robot at a distance of approximately 1.5 m. NAO was standing on a
220 desk (see Figure 1 for the setup). The transition between the two sessions was immediate with no break in
221 between. Both sessions happened in the same setting. The only change was that the researcher placed a
222 wooden box (80cm x 50cm x 50cm) on the table in front of the robot right before the LEGO session started.
223 The box was used to provide the participants with a comfortable height to complete the building tasks.
224 Together with the box, the participants were given a set of LEGO blocks (pre-built shapes) with which they
225 were to perform the tasks (see Figure 2).

226 The between-subjects design required each person participating in either one of the two conditions. In the
227 baseline condition the robot performed free from errors. In the experimental condition the robot committed
228 two social norm violations and two technical failures each in both sessions. After each robot error, the
229 researchers waited for the situation to unravel without them interfering. In many cases the participants
230 showed a reaction that confirmed that they had noticed the error (e.g., some participants laughed or frowned)
231 and then moved on. The researchers only intervened in the rare cases where the interaction was severely
232 interrupted, for example, when the participant directly addressed the researchers and commented on the
233 error. In this case, the researcher simply asked the participant to continue interacting with the robot, in
234 order to limit the interference as much as possible.

235 The three starting questions in the interview session and the first two building tasks were meant as
236 an introduction and were not varied in order. Therefore, the robot errors occurred in the randomized
237 questions/tasks only. Tables 1 and 2 give an overview on the questions and tasks and which errors occurred

¹ <https://www.ald.softbankrobotics.com/en/cool-robots/nao>



Figure 1. Study setup with the participant interacting with the robot and two researchers seated behind a bookshelf who supervised the technology

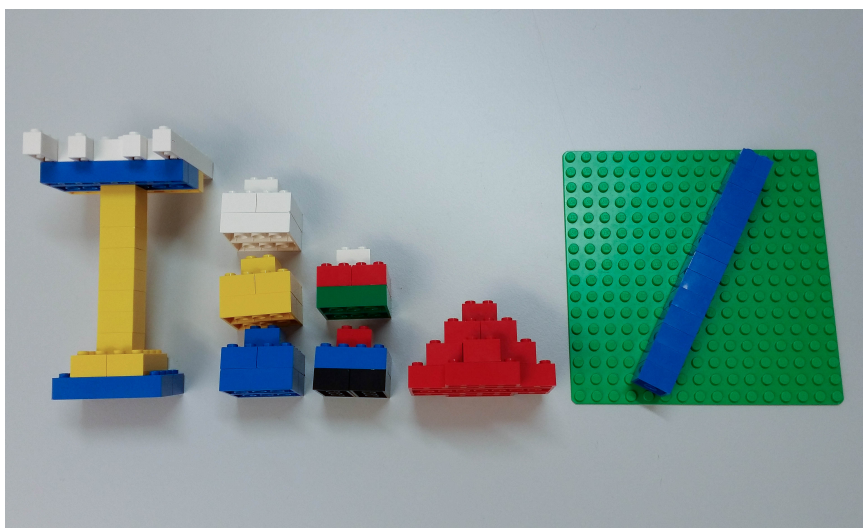


Figure 2. LEGO blocks that were provided to the participants

238 together with which question or task. The questions were similar in both conditions. The difference between
239 the baseline and the experimental condition was achieved by the presence or absence of the robot errors.

240 The induced errors were mainly modelled based on our previous findings on typical robot errors as
241 reported in Giuliani et al. (2015); Mirnig et al. (2015). Only LEGO task number 7 in the *error* condition
242 was inspired by unusual requests as reported in Salem et al. (2015).

243 The setup of our user study is based on real-life HRI. It is data-driven in representing actual error
244 situations and corresponding robot errors that occur when humans interact with state-of-the-art social
245 robots, which makes our setup ecologically valid.

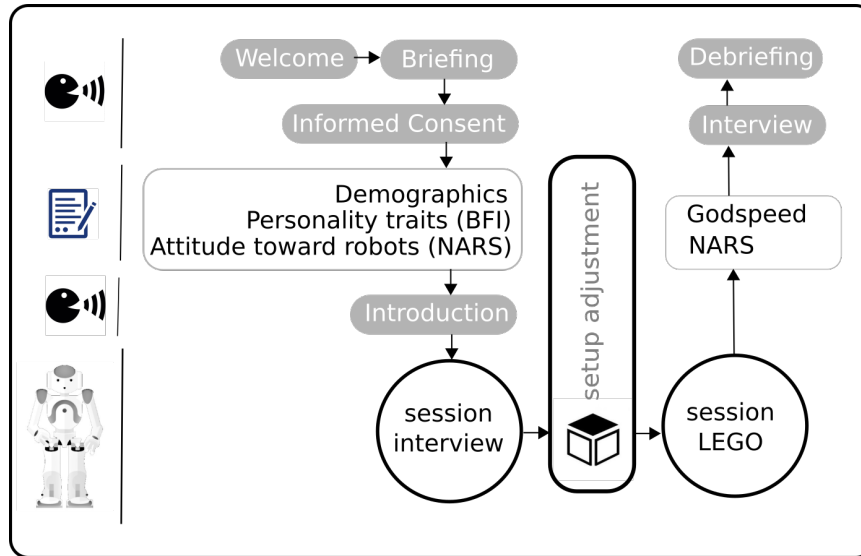


Figure 3. Study procedure

Table 1. Interview Session. The questions comprised two Social Norm Violations (SNV) and two Technical Failures (TF)

	#	Question	Error Type	Error
fixed order	1	What do you think is a robot?	-	none
	2	Which three properties come to your mind when you think about robots?	-	none
	3	Which robots do you know?	-	none
randomized order	4	Would you like a robot that assists you with household chores?	SNV	The robot waits 15 seconds until it speaks again.
	5	Why do you think some people are afraid of robots?	SNV	The robot starts speaking after 2.5 seconds, cutting off the participant.
	6	Which skills would you like for a robot to have?	-	none
	7	In which areas could humanoid robots be helpful?	-	none
	8	Have you interacted with a robot before?	TF	The robot starts speaking but cuts the sentence off after “interac”.
	9	Is hard- or software more important to you?	TF	The robot repeats the sentence 6 times.
	10	Which tasks would you never entrust a robot with?	-	none

246 **2.3 User Study Procedure**

247 The participants were welcomed to the User Experience and Interaction Experimentation Lab. After a
 248 short briefing, they were asked to sign an informed consent. Next, the participants were asked to complete
 249 questionnaires to assess their demographics, personality traits and attitude toward robots. The participants
 250 were introduced to the robot and they were given an overview on the process of the user study. As soon
 251 as the participants took their position opposite the robot, the user study began. First, the participants

Table 2. LEGO Session. The tasks comprised two Social Norm Violations (SNV) and two Technical Failures (TF).

	#	Task	Error Type	Error
fixed order	1	Place all single-color blocks on top of each other. The order does not matter [Participant performs task]. Unfortunately, the colors do not match how I imagined. Please take the blocks apart again.	-	none
	2	What animal comes to your mind? Please draw it with the blue blocks onto the green board and show it to me.	-	none
randomized order	3	Pick the multi-color block you like least. Disassemble it and build something new.	-	none
	4	Build a tower from all blocks that have red pieces in them.	-	none
	5	Build a bridge from four blocks that gets as long as possible [Participant performs task]. Wonderful! Please disassemble the bridge into the four original blocks.	-	none
	6	Count how many parts the red pyramid is made of. If you need to disassemble the pyramid to count the bricks put it back together in the end. Tell me the number.	-	none
	7	Place all single-color blocks on the right side and the remaining blocks on the left. (<i>no error condition</i>)/Throw three blocks on the floor at once! (<i>error condition</i>)	SNV	In the <i>error</i> condition, instead of giving the sorting task to the participant, the robot instructs the participant to throw three blocks on the floor at once.
	8	Place all blocks in a row sorting them by size. Begin with the smallest.	SNV	The robot waits 15 seconds until it speaks again.
	9	Build something creative from the yellow and the blue block.	TF	The robots repeats the word yellow as if stuck in a loop (“Build something creative from the yellow, yellow, yellow, ...”)
	10	Which facial expression depicts your current emotional state? Please draw the expression with the blue blocks onto the green board [Participant performs task]. Please place the picture in my hands. With the command “grasp!” I close my hands.	TF	The robot tries closing its hands but repeatedly fails to grasp the piece.

252 answered a set of questions the robot asked them (Session 1). Second, the robot instructed the participants
 253 to complete a set of building tasks with LEGO blocks (Session 2). After the interaction with the robot,
 254 the participants were again asked to complete the questionnaire assessing their attitude toward robots.
 255 They were further asked to complete a questionnaire rating the robot’s likability, anthropomorphism, and
 256 perceived intelligence. The study was finalized with a closing interview where the researcher asked the
 257 participants four open-ended questions which were followed by a short debriefing in which the purpose of
 258 the study was explained to the participants. The study procedure is depicted in Figure 3.

259 2.4 Dependent Measures

260 Before the interaction, we asked our participants to fill in the Big Five Inventory (BFI) questionnaire by
261 John et al. (2008). We used this questionnaire to analyze if people’s personality influences how they perceive
262 the robot. The BFI consists of 44 items (5-point Likert-scaled), constructing five subscales (extraversion,
263 agreeableness, conscientiousness, neuroticism, openness). This questionnaire is a well-accepted instrument
264 among psychologists to assess the personality of humans. Therefore, we chose to use it for exploring
265 potential connections between personality and how a social robot is perceived.

266 We used the Negative Attitude Towards Robots Scale (NARS), Nomura et al. (2004), to assess participants’
267 general attitude towards robots. The NARS consists of 14 items (5-point Likert-scaled) that account for
268 three scales: people’s negative attitude toward (S1) interaction with robots, (S2) social influence of robots,
269 and (S3) emotions in interaction with robots. We asked the participants to complete the questionnaire
270 before and after their interaction with the robot in order to measure if the interaction changed people’s
271 attitude. The NARS is a widely-used questionnaire in the HRI community and it provides researchers with
272 a comprehensive understanding of human fears around social robots.

273 To explore how our participants rate the robot, we used three subscales from the Godspeed Questionnaire
274 Series by Bartneck et al. (2009), i.e., anthropomorphism, likability, and perceived intelligence. Each of the
275 scales consists of five 5-point Likert-scaled items. The scales were developed in the HRI community to
276 specifically assess users’ perception of social robots. We chose the questionnaires since they are frequently
277 used and widely accepted among the HRI community. The concepts the questionnaires cover are very
278 relevant to social HRI and they represent the concepts we explore with our research. This questionnaire
279 was administered once, after our participants’ interaction with the robot.

280 2.5 Interview Data

281 We used two sources to gain qualitative data from the participants regarding their attitude toward robots.
282 First, the robot asked the participants about their opinion on robots in the interview session (see Table 1).
283 Second, in the concluding interview after the interaction and after all the other questionnaires were filled
284 in, we asked the following questions:

- 285 1. Did you notice anything special during your interaction with the robot that you would like to tell us?
- 286 2. Did your attitude toward robots change during the interaction?
- 287 3. What would you change about the interaction with the robot?
- 288 4. What did you think when the robot made a mistake? (This question was only asked for participants
289 who took part in the *error* condition).

290 2.6 Participants

291 A total of 45 participants took part in our user study (25 males and 20 females). The participants were
292 recruited over a university mailing list and social media. They were primarily university students and
293 they had not previous experience with robots. Their age ranged from 16 to 76 years, with a mean age
294 of 25.91 years ($SD = 10.82$). As regards conditions, 21 participants completed the *error* condition and
295 24 the *no error* condition. The participants’ technology affinity was rated on average with a mean of
296 3.09 ($SD = 1.49$; 5-point Likert-scaled ranging from 1 - “not technical” to 5 - “technical”) and their
297 pre-experience with robots was below average with a mean of 1.96 ($SD = 0.82$; 5-point Likert-scaled
298 ranging from 1 - “never seen” to 5 - “frequent usage”).



Figure 4. Participant interacting with the robot during the LEGO building session

299 2.7 Manipulation Check

300 In order to verify that the manipulation programmed into the robot's behavior was effective, we analyzed
301 the videos of the interactions. Out of the 21 participants of the *error* condition, 18 exhibited clearly
302 noticeable reactions upon the robot's faults (e.g., laughing, looking up from the LEGO at the robot,
303 annoyed facial expression). During the closing interview with the researcher, 15 of the 21 participants
304 stated that they noticed the robot making errors. All three persons who had not shown reactions upon the
305 robot's errors in the video mentioned them in the interview. We, therefore, conclude that our manipulation
306 was effective.

3 RESULTS

307 We used non-parametric statistical test procedures for data analysis, since our data was mostly not
308 normally distributed (Kolmogorov-Smirnov test). Mann-Whitney U tests were used to compare between
309 two independent samples (between the two conditions and between the genders). Wilcoxon rank-sum tests
310 were used to compare paired-samples (ratings of the same scales before and after the interaction).

311 We coded the qualitative data from both interviews thematically (the one the robot conducted and the
312 concluding interview after the interaction). We further annotated the video recordings from the participants'
313 interaction to investigate their social signals when experiencing an error situation with the robot. Figure 4
314 shows a participant interacting with the robot during the LEGO building session. The coding was performed
315 from one of the authors since we coded objectively visible events only.

316 3.1 Questionnaire Data

317 The gender distribution across conditions was roughly balanced. While 24 participants (15 males and 9
318 females) interacted with a flawless robot in the *no error* baseline condition, 21 participants (10 males and
319 11 females) were interviewed by an error-prone robot in the *error* experimental condition.

320 3.1.1 Participants' personality

321 We explored if our participants' personality influenced their rating of the robot by measuring five major
322 personality traits. The scales of the BFI are constructed with semantic differential items that measure the
323 participants' position between two poles (e.g., 1 - introvert to 5 - extravert). The arithmetic mean of these
324 items with no emphasis on either one of the poles is 2.5.

325 **Scale Reliability.** The subscales extraversion, neuroticism, and openness resulted in high reliability (Cron-
326 bach's $\alpha = .82, .81, .85$). The reliability for the conscientiousness scale was acceptable ($\alpha = .71$) and the
327 one for agreeableness borderline acceptable ($\alpha = .61$).

328 **Participants' overall personality.** The results showed that the participants were slightly more extroverted
329 ($mean = 3.34, SD = .72$), conscientious ($mean = 3.42, SD = .57$), and open ($mean = 3.38, SD = .79$)
330 than the arithmetic mean. They were rather agreeable ($mean = 3.79, SD = .47$), and slightly less neurotic
331 than average ($mean = 2.91, SD = .73$).

332 **Participants' personality compared between conditions.** We performed Mann-Whitney U tests to
333 explore if participants' personality profile differed between conditions. The tests for all three subsca-
334 les were non significant, showing that participants' personality profile did not differ between people who
335 completed the *error* condition and people who completed the *no error* condition ($U \geq 235, z \geq -.388, p \geq$
336 $.553, r \geq .03$).

337 3.1.2 Participants' Negative Attitude Toward Robots

338 We measured people's negative attitude toward robots for two reasons. First, we wanted to assess our
339 participants' general attitude. Therefore, we administered the NARS questionnaire before the participants'
340 interaction with the robot. Second, we assumed that participants' attitude would be affected through
341 the high number of errors. Therefore, we administered the questionnaire a second time, following the
342 interaction. The individual NARS items range from 1 - "I strongly disagree" to 5 - "I strongly agree"².
343 This means that low scale values indicate that people have a more positive attitude towards robots and high
344 scale values denote a rather negative attitude.

345 **Scale Reliability.** We checked the reliability for all three subscales, before and after the interaction. The
346 reliability for S1 before interaction resulted in borderline acceptable reliability (Cronbach's $\alpha = .64$),
347 S1 after interaction in acceptable reliability ($\alpha = .74$). The reliability for S2 before interaction was
348 too low ($\alpha = .51$). To increase reliability, we excluded item 2 (I feel that in the future society will be
349 dominated by robots), and we recalculated the scale which resulted in borderline acceptable reliability
350 ($\alpha = .62$). S2 after interaction was recalculated accordingly after excluding item 2 ($\alpha = .77$). S3 resulted in
351 borderline acceptable reliability both before and after interaction (Cronbach's α before interaction = .62,
352 after interaction = .67).

353 **Participants' overall negative attitude toward robots.** While our participants' rating for S2 and S3
354 resulted in a neutral standpoint, the rating for S1 showed that participants have a rather positive to neutral
355 attitude toward interacting with robots (mean values before interaction are presented in Table 3).

356 **Participants' negative attitude toward robots compared between before and after interaction.** We
357 were interested in investigating if our participants' negative attitude toward robots was influenced by

² Nomura et al. (2004) recommend calculating the NARS scales by summing up the item values. Since the scales are constructed of a varying number of items, the scale scores are in that case not comparable at first sight (Scale 1 would range from 6-30, Scale 2 from 5-25, Scale 3 from 3-15). Therefore, we calculated the scale values by averaging the scale items. With this, the values of the three scales become comparable more quickly and they also correlated with the range of the individual items.

Table 3. Means (SD) of the NARS questionnaire before and after the interaction (*error* and *no error* combined)

NARS Scale	before interaction	after interaction
S1: Negative Attitude toward Situations of Interaction with Robots	mean = 2.07(SD = .59)	mean = 2.09(SD = .67)
S2: Negative Attitude toward Social Influence of Robots	mean = 2.94(SD = .77)	mean = 3.11(SD = .89)
S3: Negative Attitude toward Emotions in Interaction with Robots	mean = 2.99(SD = .87)	mean = 2.79(SD = .77)

358 their interaction with the robot. We conducted Wilcoxon rank-sum tests to evaluate if the ratings differed
 359 significantly before and after the interaction. The results showed that there was no significant difference in
 360 NARS ratings before and after the interaction with the robot (S1: $W = 248.00, z = -.59, p = .558, r =$
 361 $-.06$; S2: $W = 460.00, z = 1.66, p = .097, r = -.18$; S3: $W = 234.50, z = -1.81, p = .071, r = -.19$).
 362 The means for the three scales before and after the participants' interaction with the robot are provided in
 363 Table 3.

364 **Participants' negative attitude toward robots compared between conditions.** We explored if partici-
 365 pants' rating after their interaction with the robot differed between the *error* and *no error* condition. We
 366 conducted Mann-Whitney-U tests for the scales completed after interaction. However, none of the scales
 367 resulted in significant differences between the conditions (S1: $U = 277.50, z = .85, p = .395, r = .13$;
 368 $S2 :U=324.50, z=1.66, p=.098, r=.25$; S3: $U = 277.00, z = .58, p = .564, r = .09$).

369 **Participants' negative attitude toward robots compared between the genders.** We performed Mann-
 370 Whitney-U tests to assess if the NARS ratings differed between male and female participants. The ratings
 371 for S2 and S3 (both before and after interaction) did not differ significantly. However, both ratings for
 372 S1 differed significantly between men and women (S1 before interaction: $U = 419.50, z = 3.89, p =$
 373 $.000, r = .58$; S1 after interaction: $U = 341.50, z = 2.41, p = .016, r = .36$). This result yielded in a large
 374 (before) and medium (after) effect size. For an overview on the means refer to Table 4. Even though males
 375 and females rated their potential interaction with a robot as rather positive, males ratings are significantly
 376 more positive than those of the female participants.

Table 4. NARS S1 means (SD) before and after interaction for male and female participants

NARS S1	males	females
before	mean = 1.77, SD = .54	mean = 2.46, SD = .42
after	mean = 1.87, SD = .55	mean = 2.35, SD = .73

377 3.1.3 Participants' Rating of the Robot

378 We measured how people rated the likability, anthropomorphism, and perceived intelligence of the
 379 robot after interacting with it. To do so, we used the three corresponding subscales of the Godspeed
 380 questionnaire, each of which consists of five semantic differential items. The items are constructed with
 381 semantic differential items that measure the participants' position between two poles. Therefore, the
 382 arithmetic mean of these items with no emphasis on either one of the poles is 2.5. The calculated likability
 383 score ranges from 1 - "dislike" to 5 - "like", anthropomorphism from 1 - "fake" to 5 - "natural", and
 384 perceived intelligence from 1 - "incompetent" to 5 - "competent".

385 **Scale Reliability.** The anthropomorphism and perceived intelligence scales resulted in acceptable reliability
 386 (Cronbach's $\alpha = .78, .79$), and likability in high reliability ($\alpha = .83$).

387 **Participants' overall rating of the robot.** Our participants rated the robot slightly less anthropomorphic
 388 than the arithmetic mean ($mean = 2.16, SD = .74$), more intelligent ($mean = 3.28, SD = .69$), and
 389 considerably more likeable ($mean = 4.10, SD = .63$).

390 **Participants' rating of the robot compared between conditions.** In order to explore if people who
 391 experienced erroneous robot behavior rated the robot differently from those participants who had interacted
 392 with a flawless robot, we conducted Mann-Whitney-U tests (see Table 5). While the mean ratings for
 393 anthropomorphism and perceived intelligence did not differ significantly between conditions, participants'
 394 rating of the robot's likability differed significantly between conditions. People who interacted with an
 395 erroneous robot, liked the robot significantly more than people who interacted with a flawless robot. This
 396 difference yielded in a medium effect size.

Table 5. Godspeed means (SD) compared between conditions (* denotes significant differences)

Godspeed Scale	error	no error	Mann-Whitney-U
Anthropomorphism	$mean = 1.97, SD = .66$	$mean = 2.33, SD = .78$	$U = 182.00, z = -1.60, p = .109, r = .24$
Likability*	$mean = 4.30, SD = .49$	$mean = 3.93, SD = .70$	$U = 340.00, z = 2.02, p = .044, r = .30$
Perceived Intelligence	$mean = 3.33, SD = .62$	$mean = 3.23, SD = .76$	$U = 267.50, z = .35, p = .723, r = .05$

397 **Participants' rating of the robot compared between the genders.** We conducted further Mann-Whitney-
 398 U tests to detect potential differences in robot ratings between the genders. The tests showed that none
 399 of the three scales resulted in different ratings for male and female participants (anthropomorphism:
 400 $U = 290.50, z = .93, p = .352, r = .14$; likability: $U = 317.50, z = 1.55, p = .121, r = .23$; perceived
 401 intelligence: $U = 323.00, z = 1.68, p = .094, r = .25$). We further checked if our participants' age, their
 402 pre-experience with robots, and their technological affinity influenced how the robot was rated. None of
 403 these attributes resulted in significant differences.

404 Given our results, we can infer the following for our previously postulated hypotheses. Our participants
 405 liked the robot that made errors significantly more than the flawless robot which confirms our hypothesis
 406 1. The hypotheses 2 and 3 have to be rejected since the robot committing errors did neither result in
 407 significantly higher anthropomorphism nor in significantly lower perceived intelligence ratings.

408 3.2 Qualitative Data

409 For the qualitative data analysis we annotated the video recordings of the interview and LEGO sessions
 410 from the *error* condition. We hand-coded the social signals the participants showed toward the robot, not
 411 toward the researcher and which were objectively countable. Ambiguous events were discarded. For two of
 412 the participants, there was no video data due to technical problems from the recording equipment. The
 413 video data reported, is based on the remaining 19 participants that completed the error condition. The data
 414 from the concluding interview was coded thematically in order to support our findings.

415 In this results section, we will report those findings from the qualitative data that are related to our
 416 research topic of robot errors.

417 3.2.1 Interview and LEGO Session

418 **Interview Session.** NAO began the interview with asking the participants to state their definition of a robot.
419 The majority of people provided a very technical definition: 17 people used the word machine, 10 the
420 word device, and 10 referred to a robot as some other technical object. While 2 people directly referred to
421 NAO as being a robot (“NAO, you are a robot.”), 4 participants used an “organic” noun (i.e., human, life
422 form, creature). However, they still used a technical adjective to further specify that noun (i.e., mechanical,
423 artificial, electronic, technical). Two participants provided unrelated answers.

424 We had the above question included in the robot’s questionnaire to gather people’s general standpoint on
425 robots. Since most of the participants regarded a robot as a technical object, we assumed that they would
426 want it to work reliably. In order to back our assumption up, the robot’s next question targeted the three
427 most prominent qualities people attribute with a robot. Again, many participants listed technical terms
428 ($N = 24$; e.g., mechanical, electronic, programmed). While 11 participant attributed a practical quality to
429 robots (e.g., helpful, efficient, diligent), 3 people said robots were intelligent, 6 people pointed out that
430 robots are controlled by humans (e.g., there is human intelligence in the background, not very intelligent,
431 no free will). As regards performance, 3 people referred to robots as precise/reliable, 1 participants said
432 that robots would do what they are meant to, given they are programmed correctly, and only one person
433 said that robots often make errors. This confirms our previous assumption that people assume robots to
434 perform error-free.

435 The questions reported above were asked at the beginning of the interview. In order to make the participant
436 familiar with the situation, no errors were included in here, irregardless of the condition (for a complete
437 description of the user study procedure refer to Section 2.2). Therefore, the answers were not influenced by
438 the fact that the robot made or did not make mistakes. The following questions, however, contained robot
439 errors in the *error* condition.

440 Upon asking the participants which skills they would want a robot to have, 8 participants referred to
441 robots as error-free (e.g., should do what people tell it to do, work reliably, make no mistakes). Other
442 skills included that the robot should be helpful and take on work that is too difficult/tedious/dangerous for
443 humans ($N = 13$), it should be communicative and understand the human ($N = 5$), it should be easy to
444 handle ($N = 3$), and it should be witty ($N = 2$).

445 **LEGO Session.** The robot asked the participants to express their current emotional state with LEGO bricks.
446 The emotional state declarations were classified through lip and/or eyebrow shape (for an example see
447 Figure 5). Most of the emotional state declarations were closely modelled to emoticons that are widely
448 used in social media. Depcitions that could not clearly be matched to an emotion were excluded (no data
449 entries in Figure 6). No apparent difference of participants’ emotional state could be detected between the
450 conditions. While the majority of participants was happy, only a few indicated a neutral expression. In
451 the baseline condition, one participant reported a puzzled feeling and one felt silly. In the experimental
452 condition, one participant indicated to be sad, one surprised. For an overview on all emotions refer to
453 Figure 6.

454 In the *error* condition, the robot failed to grasp the LEGO board that the participants were supposed
455 to hand over. Since the participants were instructed to tell the robot to grasp, we wanted to know how
456 often participants were willing to repeat their instructions. The number of expressed instructions (“grasp!”)
457 ranged from 2 to 7 ($mean = 4.16$, $SD = 1.21$). This result lets us assume that people are to some extent
458 patient with a faulty robot.

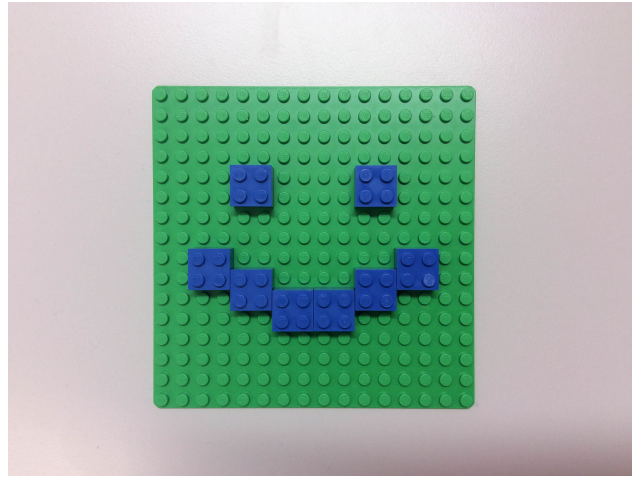


Figure 5. An example of how the participants showed their current emotion to NAO during the LEGO session

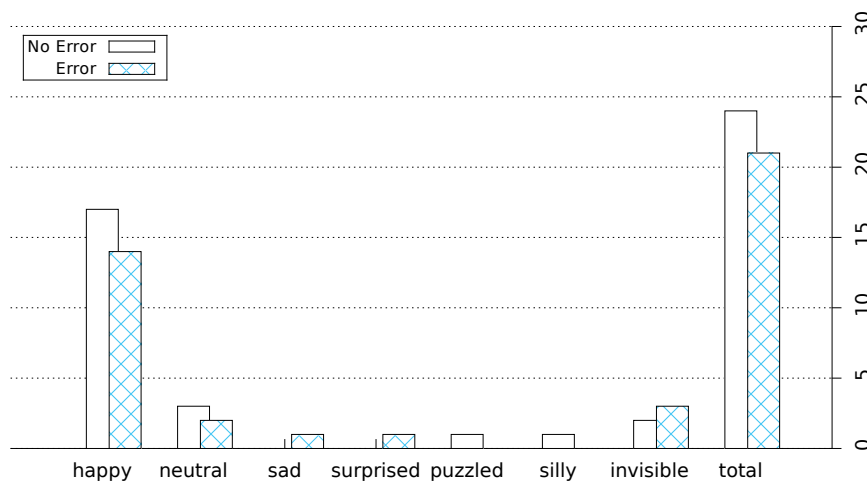


Figure 6. Emotions the participants expressed during the LEGO session

459 Upon placing an unusual request in the *error* condition, the participants' willingness to comply was
 460 striking. A total of 17 participants threw LEGO blocks to the floor when asked to do so and 2 participants
 461 bent down and placed them on the floor, but no one refused to carry out the robot's request. The fact
 462 that the participants complied with the robot's unusual request links up with the research of Salem et al.
 463 (2015). The authors report that although people seemed to know that the robot's request was not right (the
 464 researchers made the robot ask a number of unusual things of the participants, such as throwing someone's
 465 personal mail in a garbage can), people complied as long as the action was not fatal and could be undone.

466 **Social Signals.** As we intended, the participants correctly interpreted the majority of SNVs and TFs as
 467 error situations. The effectiveness manifests in the circumstance that most participants produced social
 468 signals when the robot made an error. Only the error where the robot waited for 15 seconds until it spoke
 469 was not recognized in 3 cases in the interview and in 7 cases in the LEGO session. The video footage
 470 showed that during the LEGO session, the participants were simply preoccupied with the previous task. This
 471 means that they were still dealing with the LEGO bricks (e.g., disassembling, counting, assembling, etc.)
 472 and, thus, did not pay attention to the robot's long silence. During the interview session, three participants

473 were more patient than the rest of our sample and just waited for the robot to continue. The SNV in the
 474 interview session where the robot cut the participant off, did not work in one case. This participant provided
 475 such a short but coherent answer that he was finished by the time the robot started speaking.

476 Each of the 19 participants experienced 8 error situations, which results in 152 error situations. From
 477 those, 11 were not recognized as error (see above) and in 19 cases, the participants did not show a reaction
 478 towards the robot. This leaves us with 122 error situations in which the participants showed 1 or more
 479 social signals (maximum 5). See Table 6 for an overview on the mean number of social signals per error
 480 situation.

Table 6. Mean number of social signals and standard deviation (SD) per error situation

Error Situation	Mean	SD
Interview - robot waits 15 seconds (SNV)	1.69	.946
Interview - robot cuts participant off (SNV)	1.44	.784
Interview - robot stops mid-word (TF)	.95	.911
Interview - speech loop (TF)	1.63	1.065
LEGO - throw block on the floor (SNV)	1.16	.765
LEGO - robot waits 15 seconds (SNV)	1.00	.953
LEGO - speech loop (TF)	2.00	1.106
LEGO - robot fails to grasp (TF)	2.63	1.26

481 The mean number of social signals expressed during a SNV is 1.36 ($SD = .56$) and during a TF 1.53
 482 ($SD = .72$). A Kolmogorov-Smirnov test for normality over the differences of the variable scores indicated
 483 that the data are normally distributed ($D(19) = .131, p = .200$). We performed a paired-samples t-test and
 484 found that the amount of social signals participants produced did not differ significantly between SNV and
 485 TF ($t(18) = -1.112, p = .281, d = .27$). Table 7 gives an overview on how many social signals were
 486 made for each category in each type of error situation. The table also shows which kinds of social signals
 487 were grouped in the categories. Our analysis contains only social signals that were made towards the robot.
 488 Signals towards the present experimenters were not included in our analysis (e.g., verbal statements to
 489 the experimenter, head turns in the direction of the experimenter). We hand-coded the data by counting
 490 the objectively perceivable events. Thereby, we distinguished a head tilt (head moves sideways with gaze
 491 staying in place) from a shift in gaze (the participant's gaze shifts visibly from e.g., the robot to the LEGO
 492 parts). Head turns (head movements with the gaze leaving the scene) were all directed towards the present
 493 experiment and, thus, disregarded.

494 A Kolmogorov-Smirnov test for normality over the frequency differences of the variable scores for
 495 the speech category indicated that the data deviate from normal distribution ($D(19) = .250, p = .003$).
 496 Therefore, we performed Wilcoxon signed-rank tests to assess the differences in frequencies for each
 497 category. Table 8 provides an overview on the mean number of social signal of each category per error
 498 situation type. The results show that during technical failures people made significantly more facial
 499 expressions, head movements, body movements, and gaze shifts.

500 3.2.2 Concluding Interview by the Researcher

501 After the participants finished interacting with the robot and after they completed the post-interaction
 502 questionnaires (NARS after interaction and Godspeed), they were asked four open-ended questions in the

Table 7. Overview on social signal categories and frequencies per error type

Category	Social Signals	Frequencies in SNV	Frequencies in TF
Speech	Statements, questions	13	16
Smile/laughter	Smiles, laughs, giggle	29	30
Facial expressions	Frown, raised eyebrows, corners of the mouth lowered eyes wide open	6	17
Head movements	tilted head, nodding	5	12
Body movements	lean forward, step back, touch face, adjust glasses put hands on hip, put hands behind back, take hands out from pockets raise arm and dance, sway, snap fingers, move LEGO parts around in front of the robot	8	19
Gaze shift	shift gaze to or away from robot, wandering gaze	26	43
Total number of social signals		87	137

Table 8. Social signals shown during social norm violations and technical failures

Social Signal	Social Norm Violation	Technical Failure	Wilcoxon signed-rank		
	Mean (SD)	Mean (SD)	Z	p-value	r-value
Speech	0.68 (0.820)	0.84 (0.958)	0.758	0.448	0.12
Smile/laughter	1.53 (1.219)	1.58 (0.902)	-0.074	0.941	-0.01
Facial expressions	0.32 (0.582)	0.89 (0.809)	-2.147	0.032	-0.35
Head movements	0.26 (0.562)	0.63 (1.165)	-2.121	0.034	-0.34
Body movements	0.42 (0.607)	1.00 (0.816)	-2.484	0.013	0.40
Gaze shift	1.37 (0.831)	2.26 (1.098)	-3.090	0.002	0.50

503 final interview. While the questions 1-3 asked about some general aspects of the participants' impression
 504 of the interaction and the robot, question 4 specifically targeted the robot's errors (see Section 2.5 for
 505 the specific questions). Therefore, question 4 was only asked for participants in the *error* condition. The
 506 resulting data was analyzed through an affinity diagram (Holtzblatt et al. (2004)). An affinity diagram is a
 507 method for organizing ideas, challenges, and solutions into a wall-sized hierarchical diagram.

508 In question 1, participants were asked to report anything particular they had noticed during their interaction
 509 with the robot. Here, 12 participants reported that the robot had made some mistakes (e.g., *it went in a loop*;
 510 *it cut my word*). The participants' answers to question 2 did not include any mentions about the robot's
 511 mistakes. In question 3, seven participants reported that they would like to change the faulty robot behavior
 512 (e.g., *fix the technical bugs*; *it does not leave time for you to respond*; *loops*).

513 With the final question in the interview, we specifically targeted the robot's errors, in asking what the
 514 participants thought of the robot making mistakes. While 7 participants uttered specifically negative aspects
 515 (e.g., *unpleasant*; *confusing*; *that's just what one would expect from technology*; *I was unsure if the*
 516 *interaction had stopped*; *I thought I had made a mistake*), 10 participants uttered positive feelings when
 517 asked about the fact that the robot made mistakes (e.g., *funny*; *friendly*; *it was great that the robot did not*
 518 *make it look like I made a mistake*; *I don't like it less because of the mistakes*; *it would be scary if all went*
 519 *smooth because that would be too human-like*).

4 DISCUSSION

520 Our results showed that the participants liked the faulty robot significantly more than the flawless one. This
521 finding confirms the *Pratfall Effect* which states that people's attractiveness increases when they make
522 a mistake as shown by Aronson et al. (1966). Therefore, the psychological concept can successfully be
523 transferred from interpersonal interaction to HRI. Upon the attempt of including socially acting robots into
524 this concept, we can extend it to: "*Imperfections and mistakes carry the potential of increasing the likability*
525 *of any social actor (human or robotic).*" The same effect was previously researched by Salem et al. (2013)
526 where incongruent behavior of a robot can be seen as a social norm violation as such behavior violates
527 participants' expectations from a *social script*. To overcome this error situation, participants changed their
528 social signals, but on the other hand they rated the likability of the robot higher. Similarly, Ragni et al.
529 (2016) showed that the participants in their study enjoyed the interaction with the faulty robot significantly
530 more, than the participants who had interacted with a flawless robot. On the other hand, their participants
531 who had interacted with the faulty robot, rated it less intelligent, less competent, and less superior, which
532 again confirms the *Pratfall Effect*.

533 The repeated evidence of this phenomenon existing in HRI, strengthens our argument to create robots
534 that do not lead to believe they perform free from errors. We recommend that robot creators design social
535 robots with their potential imperfections in mind. We see two sources for these imperfections that link
536 back to the two error types found in HRI. On one hand, creators of social robots should follow the notions
537 of interpersonal interaction to meet the expectations humans have about social actors and with it socially
538 interacting robots. On the other hand, it is advisable to embrace the imperfections of robot technology.
539 Technology that is created with potential shortcomings in mind, can be designed to include methods for
540 error recovery. Therefore, one way to go here would be to make robots understand they made an error by
541 correctly interpreting the human's social signals and indicate their understanding to the human user. Both
542 of these sources of imperfections will lead to more believable robot characters and more natural interaction.
543 Of course, this applies to social robots operating in non-critical environments. Safety-relevant applications
544 and scenarios must under all circumstances operate at zero-defect level.

545 Interestingly, we could not find a comparable effect for anthropomorphism in our data. The robot's
546 anthropomorphism level was rated similar, irregardless of the fact if the robot made errors or not. Our
547 result is different from the findings of Salem et al. (2013), who also used a human-like robot, and where
548 the participants rated the faulty robot more anthropomorphic as the flawless one. The researchers used
549 co-verbal gestures, while we programmed the robot to provide mostly random gestures to make it appear
550 more life-like. This might have in general diminished the effect of anthropomorphism in our setup (which
551 is indicated by the low overall anthropomorphism level). However, more research is required to further
552 explore the role of anthropomorphism in faulty robot behavior.

553 Contrary to our assumption, the faulty robot was not rated as less intelligent than the flawless one. This
554 seems striking since the robot made several errors over a relatively short interaction time. Furthermore,
555 most participants had noticed the robot making errors, while, at the same time, they had indicated to regard
556 a robot as something very technical that should perform reliably. One potential explanation could be the
557 fact that the induced errors were non task-related. Follow-up research is required to further explore the
558 perceived intelligence of erroneous robot behavior.

559 Upon asking the participants about their current emotional state, the majority of participants showed
560 the robot that they were happy. The participants were also quite patient and tried handing the object
561 several times, when the robot failed grasping it. All of these observations point towards the notion that

562 a faulty social robot is a more natural social robot. In our future research on this topic we will extend
563 our approach to include more user experience measures to get a more profound understanding on the
564 users' perception of the robot. For example, it will be interesting to further investigate possible impacts on
565 subjective performance and acceptance.

566 Our data showed that when people interacted with a social robot that made an error, they were likely to
567 show social signals in response to that error. In our previous research we performed an analysis of video
568 material in which robot errors occurred unintentionally and we found that users showed social signals in
569 about half the interactions, Mirnig et al. (2015). In the herein reported study, however, most participants
570 showed at least one social signal per error situation. We explain this difference in part with the high error
571 rate (8 errors in an average total interaction time of about 12 minutes). Users seem to anticipate the robot
572 making more errors once they experienced it is not flawless and responded more frequently with social
573 signals. The reason for the increased number of social signals could also be based on the size of the robot.
574 While the majority of interactions from the previous study were with a human-sized robot at eye level, the
575 robot in our case was small and placed slightly below participants' eye level. This aspect remains to be
576 studied further.

577 With our results we show again that humans respond to a robot's error with social signals. Therefore,
578 recognizing social signals might help a robot to understand that an error happened. According to the
579 frequencies of occurrence, gaze shifts and smile/laughter carry most potential for error detection, which
580 is in line with our previous findings in Giuliani et al. (2015). Upon a detailed analysis on the categories
581 of social signals we found that people make significantly more gaze shifts during technical failures. This
582 results is in contrast to our previous findings where significantly more gaze shifts were made during social
583 norm violations. We take from this that gaze shifts are a potential indicator for robot errors, but it remains
584 to be studied if they can be used to distinguish between the two error types.

585 We also found that people made significantly more facial expressions, head- , and body movements
586 during technical failures. The increase in social signals during technical failures may be rooted in the
587 circumstance that the technical failures were more obvious in the present user study. For example, in the
588 video material from the previous study the robot failed to grasp an object that was placed in front of it. In
589 our setup, the robot failed to grasp an object that the participant handed to it, which made the participant
590 more actively perceive the robot's error.

591 Contrary to our previous findings, we did not detect significant differences in spoken social signals. This
592 could be grounded in the fact that due the setup, the robot had in general a much larger share in spoken
593 utterances.

594 In response to the robot's unusual request, most users showed social signals. The kind of signals (gaze
595 shifts and laughter) displayed the users' slight discomfort and provided evidence that they knew the robot's
596 request implied a deviation from the social script of the situation. However, most users nevertheless
597 followed the robot's order and threw the LEGO blocks to the floor. In addition to the previous results as
598 reported in Mirnig et al. (2015), this result provides further evidence that users show specific social signals
599 in response to erroneous robot errors. Future research should be targeted at making a robot understand the
600 signals and make sense of them. A robot that can understand its human interaction partner's social signals,
601 will be a better interaction partner itself and the overall user experience will improve.

602 Since most of our participants had not interacted with a robot before, a potential novelty denotes a certain
603 limitation to our results. Some participants were probably captivated with the technology, which made
604 them remain patient. It remains to be studied how such novelty wears off over time and how this influences

605 people's willingness to interact. It will, furthermore, be interesting to assess the dimensions of faults. That
606 is, how extensive can an error become until it becomes a deal-breaker. Ragni et al. (2016) already provided
607 evidence that erroneous robot behavior decreases performance of a human interacting with the robot. It
608 could also be interesting to explore how users react in case of the robot giving ambiguous information.
609 Further aspects of robot errors that are worthwhile exploring are, for example, the following. What kinds
610 of errors are forgivable and which ones are not? What is the threshold for error rate or number of errors
611 until the participants' patience is over or performance drops considerably? A lot more specific research is
612 required to understand and make use of the effects of errors in social HRI.

5 CONCLUSIONS

613 With our user study we explored how people rated a robot making errors in comparison to a perfectly
614 performing robot. We measured the robot's likability, anthropomorphism, and perceived intelligence. We
615 found that the faulty robot was rated as more likeable, but neither more anthropomorphic nor less intelligent.
616 We recommend robots to be designed with their possible shortcomings in mind as we believe that this will
617 result in more likeable social robots. Similar to interpersonal interaction, imperfections might even have a
618 positive influence in terms of likability. We expect social HRI that embraces the imperfectness of today's
619 robots to result in more natural interaction and more believable robot characters.

620 Our results confirm existing HRI research on robot likability such as Salem et al. (2013) and Ragni et al.
621 (2016), hinting at error-prone robots supposedly resulting in more believable robots. Our work successfully
622 proves the existence of the psychological concept *Pratfall Effect* in HRI and suggests that it should be our
623 community's aim to bear potential shortcomings of social robots in mind when creating them. The nature
624 and extent of errors that can be handled through the interactional design remains yet to be studied.

625 With our results we could again show that humans respond to faulty robot behavior with social signals. A
626 robot that can recognize these social signals can, in subsequence, understand that an error happened. We
627 detected gaze shifts and laughter/smiling as the most frequently shown social signals, which is in line with
628 our previous research.

629 We see the following next steps to the ambitious goal of creating social robots that are able to overcome
630 an error situation. First, it needs to be studied how we can let robots understand that an error occurred.
631 Second, robots must be enabled to communicate about such errors. Third, robots need to know how to
632 behave in an error situation in order to effectively apply error recovery strategies.

FUNDING

633 We gratefully acknowledge the financial support by the Austrian Federal Ministry of Economy, Family
634 and Youth and the National Foundation for Research, Technology and Development (Christian Doppler
635 Laboratory for "Contextual Interfaces"). This work was additionally funded in part by the European
636 Commission in the project ReMeDi (Grant No. 610902).

ACKNOWLEDGMENTS

637 The authors of this paper would like to thank Michael Miksch for his contribution in performing the user
638 study.

REFERENCES

- 639 Aronson, E., Willerman, B., and Floyd, J. (1966). The effect of a pratfall on increasing interpersonal
640 attractiveness. *Psychonomic Science* 4, 227–228
- 641 Bartneck, C., Kulić, D., Croft, E., and Zoghbi, S. (2009). Measurement instruments for the anthropomor-
642 phism, animacy, likeability, perceived intelligence, and perceived safety of robots. *International journal*
643 *of social robotics* 1, 71–81
- 644 Brooks, D. J., Begum, M., and Yanco, H. A. (2016). Analysis of reactions towards failures and recovery
645 strategies for autonomous robots. In *Proceedings of the IEEE International Symposium on Robot and*
646 *Human Interactive Communication (RO-MAN 2016)* (IEEE), 487–492
- 647 Bruckenberger, U., Weiss, A., Mirnig, N., Strasser, E., Stadler, S., and Tscheligi, M. (2013). The good, the
648 bad, the weird: Audience evaluation of a “real” robot in relation to science fiction and mass media. In
649 *Proceedings of the International Conference on Social Robotics* (Springer), 301–310
- 650 Gehle, R., Pitsch, K., Dankert, T., and Wrede, S. (2015). Trouble-based group dynamics in real-world
651 hri—reactions on unexpected next moves of a museum guide robot. In *Proceedings of the International*
652 *Symposium on Robot and Human Interactive Communication* (IEEE), 407–412
- 653 Giuliani, M., Mirnig, N., Stollnberger, G., Stadler, S., Buchner, R., and Tscheligi, M. (2015). Systematic
654 analysis of video data from different human-robot interaction studies: A categorisation of social signals
655 during error situations. *Frontiers in Psychology* 6
- 656 Gompei, T. and Umemuro, H. (2015). A robot’s slip of the tongue: Effect of speech error on the familiarity
657 of a humanoid robot. In *Proceedings of the International Symposium on Robot and Human Interactive*
658 *Communication* (IEEE), 331–336
- 659 Hayes, C. J., Maryam, M., and Riek, L. D. (2016). Exploring implicit human responses to robot mistakes
660 in a learning from demonstration task. In *Proceedings of the International Symposium on Robot and*
661 *Human Interactive Communication* (IEEE), 246–252
- 662 Holtzblatt, K., Wendell, J. B., and Wood, S. (2004). *Rapid contextual design: a how-to guide to key*
663 *techniques for user-centered design* (Elsevier)
- 664 John, O. P., Naumann, L. P., and Soto, C. J. (2008). Paradigm shift to the integrative big five trait taxonomy.
665 *Handbook of personality: Theory and research* 3, 114–158
- 666 Knepper, R. A., Tellex, S., Li, A., Roy, N., and Rus, D. (2015). Recovering from failure by asking for help.
667 *Autonomous Robots* 39, 347–362
- 668 Lee, M. K., Kielser, S., Forlizzi, J., Srinivasa, S., and Rybski, P. (2010). Gracefully mitigating breakdowns
669 in robotic services. In *Proceedings of the 5th ACM/IEEE international conference on Human-robot*
670 *interaction* (IEEE Press), 203–210
- 671 Lohse, M. (2011). The role of expectations and situations in human-robot interaction. *New Frontiers in*
672 *Human-Robot Interaction* , 35–56
- 673 Mirnig, N., Giuliani, M., Stollnberger, G., Stadler, S., Buchner, R., and Tscheligi, M. (2015). Impact
674 of robot actions on social signals and reaction times in hri error situations. In *In Proceedings of the*
675 *International Conference on Social Robotics* (Springer), 461–471
- 676 Nomura, T., Kanda, T., Suzuki, T., and Kato, K. (2004). Psychology in human-robot communication: An
677 attempt through investigation of negative attitudes and anxiety toward robots. In *Proceedings of the*
678 *International Symposium on Robot and Human Interactive Communication* (IEEE), 35–40
- 679 Ragni, M., Rudenko, A., Kuhnert, B., and Arras, K. O. (2016). Errare humanum est: Erroneous robots in
680 human-robot interaction. In *Proceedings of the IEEE International Symposium on Robot and Human*
681 *Interactive Communication (RO-MAN 2016)* (IEEE), 501–506

- 682 Robinette, P., Wagner, A. R., and Howard, A. M. (2015). The effect of robot performance on human–robot
683 trust in time–critical situations
- 684 Salem, M., Eyssel, F., Rohlfing, K., Kopp, S., and Joublin, F. (2013). To err is human (-like): Effects of
685 robot gesture on perceived anthropomorphism and likability. *International Journal of Social Robotics* 5,
686 313–323
- 687 Salem, M., Lakatos, G., Amirabdollahian, F., and Dautenhahn, K. (2015). Would you trust a (faulty) robot?:
688 Effects of error, task type and personality on human-robot cooperation and trust. In *Proceedings of the*
689 *Tenth Annual ACM/IEEE International Conference on Human-Robot Interaction (ACM)*, 141–148
- 690 Short, E., Hart, J., Vu, M., and Scassellati, B. (2010). No fair!!: An interaction with a cheating robot. In
691 *Proceedings of the 5th ACM/IEEE International Conference on Human-robot Interaction* (Piscataway,
692 NJ, USA: IEEE Press), HRI '10, 219–226
- 693 Vinciarelli, A., Pantic, M., and Bourlard, H. (2009). Social signal processing: Survey of an emerging
694 domain. *Image and Vision Computing* 27, 1743–1759