

# Working with Troubles and Failures in Conversation between Humans and Robots: Workshop Report

Frank Förster<sup>1,\*</sup>, Marta Romeo<sup>2,3</sup>, Patrick Holthaus<sup>1</sup>, Luke Wood<sup>1</sup>, Christian Dondrup<sup>3</sup>, Joel E. Fischer<sup>10</sup>, Farhana Ferdousi Liza<sup>4</sup>, Sara Kaszuba<sup>5</sup>, Julian Hough<sup>6</sup>, Birthe Nettet<sup>3</sup>, Daniel Hernández García<sup>3</sup>, Dimosthenis Kontogiorgos<sup>7,8</sup>, Jennifer Williams<sup>9</sup>, Elif Ecem Özkan<sup>10</sup>, Pepita Barnard<sup>11</sup>, Gustavo Berumen<sup>12</sup>, Dominic Price<sup>11</sup>, Sue Cobb<sup>11</sup>, Martina Wiltschko<sup>12</sup>, Lucien Tisserand<sup>13</sup>, Martin Porcheron<sup>11,14</sup>, Manuel Giuliani<sup>16</sup>, Gabriel Skantze<sup>14</sup>, Patrick G.T. Healey<sup>10</sup>, Ioannis Papaioannou<sup>15</sup>, Dimitra Gkatzia<sup>17</sup>, Saul Albert<sup>18</sup>, Guanyu Huang<sup>19</sup>, Vladislav Maraev<sup>20</sup>, Epaminondas Kapetanios<sup>1</sup>.

<sup>1</sup>Department of Computer Science, School of Physics, Engineering and Computer Science, University of Hertfordshire, Hatfield, UK

<sup>2</sup>Department of Computer Science, The University of Manchester, Manchester, UK

<sup>3</sup>School of Mathematical and Computer Sciences, Heriot-Watt University, Edinburgh, UK

<sup>4</sup>School of Computing Sciences, University of East Anglia, Norwich, UK

<sup>5</sup>Department of Computer, Control and Management Engineering “Antonio Ruberti”, Sapienza University of Rome, Rome, Italy

<sup>6</sup>School of Mathematics and Computer Science, Swansea University, Swansea, UK

<sup>7</sup>Department of Computer Science, Humboldt University of Berlin, Berlin, Germany

<sup>8</sup>Science of Intelligence, Research Cluster of Excellence, Berlin, Germany

<sup>9</sup>School of Electronics and Computer Science, University of Southampton, Southampton, UK

<sup>10</sup>School of Electronic Engineering and Computer Science, Queen Mary University of London, London, UK

<sup>11</sup>School of Computer Science, University of Nottingham, Nottingham, UK

<sup>12</sup>ICREA, Universitat Pompeu Fabra, Barcelona, Spain

<sup>13</sup>UMR 5191 ICAR, CNRS, Labex ASLAN, ENS de Lyon, Lyon, FR

<sup>14</sup>KTH Speech Music and Hearing, Stockholm, SE

<sup>15</sup>Alana AI, London, UK

<sup>16</sup>University of the West of England Bristol

<sup>17</sup>Edinburgh Napier University, UK

<sup>18</sup>School of Social Sciences and Humanities, Loughborough University, UK

<sup>19</sup>Department of Computer Science, The University of Sheffield, Sheffield, UK

<sup>20</sup>Department of Applied IT, Univeristy of Gothenburg, Göteborg, Sweden

Correspondence\*:

Frank Förster

f.foerster@herts.ac.uk

## 2 ABSTRACT

3 This paper summarizes the structure and findings from the first *Workshop on Troubles and*  
4 *Failures in Conversations between Humans and Robots*. The workshop was organized to  
5 bring together a small, interdisciplinary group of researchers working on miscommunication  
6 from two complementary perspectives. One group of technology-oriented researchers was  
7 made up of roboticists, Human-Robot Interaction (HRI) researchers and dialogue system  
8 experts. The second group involved experts from conversation analysis, cognitive science,  
9 and linguistics. Uniting both groups of researchers is the belief that communication failures  
10 between humans and machines need to be taken seriously and that a systematic analysis  
11 of such failures may open fruitful avenues in research beyond current practices to improve  
12 such systems, including both speech-centric and multimodal interfaces. This workshop  
13 represents a starting point for this endeavour. The aim of the workshop was threefold:  
14 Firstly, to establish an interdisciplinary network of researchers that share a common interest  
15 in investigating communicative failures with a particular view towards robotic speech  
16 interfaces; secondly, to gain a partial overview of the “failure landscape” as experienced  
17 by roboticists and HRI researchers; and thirdly, to determine the potential for creating a  
18 robotic benchmark scenario for testing future speech interfaces with respect to the identified  
19 failures. The present article summarizes both the “failure landscape” surveyed during the  
20 workshop as well as the outcomes of the attempt to define a benchmark scenario.

21 **Keywords:** human-robot interaction, speech interfaces, dialogue systems, multi-modal interaction, communicative failure,  
22 repair

## 1 INTRODUCTION

23 Speech interfaces, user interfaces that allow interaction with technology through spoken commands  
24 or queries, are commonplace in many types of robots and robotic applications. Despite the progress  
25 in speech recognition and many other areas of natural language processing in recent years, failures of  
26 speech interfaces in robotic scenarios are numerous, especially in real-world situations (Porcheron  
27 et al., 2018; Fischer et al., 2019). In contrast to the common experience of failure of speech interfaces  
28 in robotics, the literature is positively skewed towards the success and good performance of these.  
29 While Marge et al. (2022) identified key scientific and engineering advances needed to enable  
30 effective spoken language interaction with robotics; little attention was given to communicative  
31 failures. To our knowledge, the documentation of failure in speech interfaces and systematic studies  
32 of such failures and their causes is exceedingly rare. Honig and Oron-Gilad (2018) provides the  
33 most in-depth literature review of prior failure-related HRI studies. The authors found that research  
34 in HRI has focused mostly on technical failures, with few studies focusing on human errors, many  
35 of which are likely to fall under the umbrella of conversational failures. In addition to this focus on  
36 technical errors, the majority of failure-related studies in HRI take place in controlled experimental

37 conditions, where ‘failures’ are explicitly designed and occur only at specific moments (Ragni  
38 et al., 2016; Washburn et al., 2020a; Cuadra et al., 2021; Green et al., 2022), instead of a natural  
39 occurrence of the interactions between humans and robots. Closer to the topic of the workshop is  
40 the recently proposed taxonomy of Tian and Oviatt (2021) that focuses on social errors in HRI and  
41 their relationship with the perceived socio-affective competence of a robot. However, while there is  
42 significant overlap between social errors, as categorized by Tian and Oviatt, and the workshop topic  
43 of conversational failure, the perspective on the role of these errors and failures in interaction as  
44 well as the view as to whether these could be overcome eventually differs significantly. While social  
45 errors should ultimately be reduced by increasing a robot’s perceived socio-affective competence, it  
46 appears unlikely that conversational failure could be totally extinguished by means of technological  
47 progress. Too frequent is their occurrence in human-human conversation and too deeply ingrained  
48 are the related repair mechanisms in the fabric of human communication.

49 To the best of our knowledge, there are currently no survey papers specifically on conversational  
50 failures in human-robot interaction, a fact that illustrates an important gap in the research landscape.  
51 To address this gap, we conducted a two-phase workshop with experts in adjacent fields. This paper  
52 presents the findings from this workshop series that brought together a multidisciplinary group of  
53 researchers from fields such as robotics, human-robot interaction (HRI), natural language processing  
54 (NLP), conversation analysis, linguistics and pragmatics. The workshop provided a platform to  
55 discuss the multitude of failures of speech interfaces openly and to point out fruitful directions for  
56 overcoming these failures systematically. The workshop focused mainly on human-robot joint action  
57 scenarios involving multimodal coordination between humans and robots, as these are the norm in  
58 scenarios where robotic speech interfaces are deployed. The identified types of failures range from  
59 failures of speech recognition to pragmatic failures and infelicities.

60 We begin by describing the aims, structure, and materials used in the workshop in Sect. 2. We then  
61 present findings that result from the workshop, including participant contributions and outcomes of  
62 the structured discussion in Sect. 3. This leads to Sect. 4, where we reflect on problems and identify  
63 themes that emerged from the workshop’s discussions before concluding the paper.

## 2 MATERIALS AND METHODS

64 The *Working with Troubles and Failures (WTF) in Conversations between Humans and Robots*  
65 workshop included a virtual gathering over two consecutive days in June 2022 and an in-person  
66 full-day meeting at the University of Hertfordshire in September 2022. Here, we sketch the structure  
67 and summarize the findings for each of these parts.

## 68 2.1 Before the Workshop

69 In order to attract workshop participants interested in an open discussion of their experience and  
70 investigations of failing speech interfaces, we directly contacted some of the potentially interested  
71 research groups within the United Kingdom. Additionally, the workshop was advertised via mailing  
72 lists relevant to the HRI (e.g. *hri-announcement*, *robotics-worldwide*, *euRobotics-dist*), natural  
73 language processing (NLP, e.g. *ACM sigsem*), and artificial intelligence communities (e.g. *ACM*  
74 *sigai-announce*). To verify participants' genuine interest in the topic and to collate information on  
75 the different types of conversational failures experienced by them, they were asked to submit the  
76 following pieces of information:

- 77 1. the number of years of experience using or developing speech interfaces,
- 78 2. an indication of what they perceive to be the most pressing issue or the biggest source of failure  
79 for speech interfaces,
- 80 3. their most memorable WTF moment, that is, which of their experiences of failure with a speech  
81 interface they remembered most vividly,
- 82 4. a summary of their motivation to attend the workshop,
- 83 5. a suggestion for a future benchmark scenario that would expose the kind of failure described in  
84 their WTF moment.

85 Applicants that stated a meaningful entry for item 4, and made some attempt to answer the other  
86 questions, were admitted to the workshop. As a result, 15 participants were admitted and initially  
87 attended the virtual part. Of these fifteen participants, eight would go on to attend the face-to-  
88 face part of the workshop. The face-to-face workshop was re-advertised via the above-mentioned  
89 mailing lists and the same set of questions and answers was used to filter out additional prospective  
90 participants. Ultimately, six new participants joined the face-to-face part of the workshop, resulting  
91 in fourteen non-speaker, non-organiser participants. Two of these attended the face-to-face workshop  
92 virtually, as we decided to go for a hybrid format in order not to exclude anyone who was not able  
93 or willing to travel on site.

94 Keynote speakers for both parts of the workshop were chosen based on their expertise in the  
95 subject area. The subject areas considered most relevant to the workshop were robotics-centred NLP  
96 on the one hand and Conversation Analysis (CA) on the other. The emphasis on CA was based on  
97 the fact that the documentation and analysis of conversational failure have been an integral part  
98 of this discipline since its very inception. Moreover, it was hoped that having keynote speakers  
99 and participants from both areas would soften discipline-specific boundaries and limitations and  
100 potentially open up new directions for future research.

### 101 2.1.1 Motivations for Attending the Workshop

102 The following is a summary of the participants' motivation for attending the workshop as extracted  
103 from the application forms:

104 Several PhD students were hoping to connect and network with other researchers working in speech  
105 interaction technologies. Multiple other researchers working on the CA-HRI interface wanted to  
106 learn more about how conversational trouble emerges, while others occupied with developing speech  
107 interfaces, or with integrating these into robots were interested in gaining a deeper understanding of  
108 current issues. Many of them were also interested in sharing their experiences with peers.

109 One researcher working in animal communication hoped to learn something from a different domain  
110 of "inter-being communication", while yet another researcher working on speech privacy wanted  
111 to connect to other researchers working on speech interfaces. One participant saw value in the aim  
112 of identifying or creating a benchmark scenario that would be able to tease out the most common  
113 failures, if they occurred - an aim explicitly set out by the workshop.

114 Another motivation of multiple participants to attend the workshop was their shared belief that a  
115 deeper analysis of communicative failures would not only help to improve future speech interfaces  
116 but also gain a deeper understanding of (human) conversations themselves.

117 Finally, a researcher interested in explainable AI was interested to see what other types of failures,  
118 apart from faulty explanations, there are and how these may connect to research in explainable AI.

## 119 2.2 Virtual Workshop

120 To facilitate participation in the virtual session of the workshop, it was divided into two half-day  
121 events. On the first day, the workshop opened with a keynote talk by Prof. Patrick Healey, Professor  
122 of Human Interaction and Head of the Cognitive Science Research Group in the School of Electronic  
123 Engineering and Computer Science at Queen Mary University of London, on "Running repairs:  
124 Coordinating meaning in dialogue" (Section 3.1.1). This was followed by participants' lightning  
125 talks on their most memorable WTF moments when working with communication between humans  
126 and robots (Section 3.2). Following the lightning talks, and based on the underlying themes identified  
127 by the organisers, participants were divided between 4 breakout rooms to continue discussing the  
128 issues they brought to the workshop. The four identified themes were: (i) Context Understanding,  
129 (ii) Handling Miscommunication, (iii) Interaction Problems, and (iv) General Failures.

130 The second day of the virtual workshop saw Dr. Saul Albert, Lecturer in Social Science (Social  
131 Psychology) in Communication and Media at Loughborough University, give a keynote talk on  
132 "Repair, recruitment, and (virtual) agency in a smart homecare setting" (Section 3.1.2). Following  
133 the talk, each group from the breakout rooms of the first day reported what was discussed and each  
134 debate was opened to all participants. The workshop ended with a short summary of the day.

## 135 2.3 Face-to-Face Workshop

136 The in-person part of the workshop was held at the University of Hertfordshire three months  
137 after the virtual event. During this full-day meeting, keynote talks were given by Prof. Gabriel  
138 Skantze, Professor in Speech Technology at KTH Royal Institute of Technology on “Building  
139 Common Ground in Human-Robot Interaction” (Section 3.1.3) and by Dr. Ioannis Papaioannou,  
140 Chief Technology Officer & Co-Founder of Alana<sup>1</sup> on “Tackling the Challenges of Open-Domain  
141 Conversational AI Systems” (Section 3.1.4).

142 Since the registration to the face-to-face workshop was also opened to participants who did not  
143 take part in the virtual workshop, new attendees were given the opportunity to present their own  
144 lightning talks on their WTF moments (Section 3.2).

145 A central part of the face-to-face workshop was the World Café session<sup>2</sup>, which provided  
146 participants an opportunity to freely discuss troubles and failures in small groups across several  
147 table topics. Based on the participants’ submitted WTF moments, and the themes from the breakout  
148 rooms of the virtual part, four themes were chosen for this session: (i) Context Understanding, (ii)  
149 Interaction Problems, (iii) Handling Miscommunication, and (iv) Suggested Benchmark Scenarios.  
150 Each theme was allocated to one table, and each table had one designated organizer. Participants  
151 and speakers were split into four different groups and moved between the tables within time slots  
152 of approximately 15 minutes per theme. The tasks of a table’s organizer were to summarize the  
153 findings and discussions from previous groups to a newly arriving group, to encourage discussions  
154 around the table topic, and to either encourage note taking or take notes themselves on a large flip  
155 chart that was allocated to each table.

## 3 RESULTS

156 In this section, we present findings from both the virtual and the face-to-face parts of the workshop,  
157 describing how the keynotes shaped the discussion and how the participant lightning talks contributed  
158 to identify some of the most pressing problems in conversations between humans and robots. Most  
159 importantly, we will present the outcomes of the structured discussion, summarising the workshop  
160 findings.

### 161 3.1 Keynotes

162 To frame the discussion on troubles and failures with experiences from different perspectives, we  
163 invited four keynote speakers from scientific areas that are concerned with research problems around

---

<sup>1</sup> <https://alanaai.com/>

<sup>2</sup> <https://theworldcafe.com/key-concepts-resources/world-cafe-method/>

164 conversations between humans and robots. This section summarises their presentations in the context  
165 of the workshop goals to scope and identify common troubles and failures in conversation between  
166 humans and robots. In the virtual part of the workshop, the first keynote (Sect. 3.1.1) provided a  
167 conversation analytical perspective on repairs and meaning in dialogue, while the second one looked  
168 at repairs but from a more applied perspective in a user's home (Sect. 3.1.2). The in-person workshop  
169 provided insights considering human-robot interactions (Sect. 3.1.3) and an industry viewpoint  
170 (Sect. 3.1.4).

### 171 3.1.1 Running Repairs: Coordinating Meaning in Dialogue

172 Healey presented the Running Repairs Hypothesis (Healey et al., 2018b), which captures the idea  
173 that successful communication depends on being able to detect and adjust to misunderstandings on  
174 the fly. The basic assumption is that no two people ever understand exactly the same thing by the  
175 same word or gesture and, as a result, misunderstandings are ubiquitous. Data from conversations  
176 support this assumption. For example, the utterance "huh?" occurs around once every 84 seconds in  
177 conversation and appears to be universal across human languages (Enfield, 2017; Dingemanse et al.,  
178 2015). Around a third of turns in ordinary conversation involve some sort of real-time adjustments  
179 in language use (Colman and Healey, 2011).

180 The processes for detecting and resolving problems with understanding have conventionally been  
181 regarded as 'noise in the signal' by the cognitive sciences (Healey et al., 2018a). However, there  
182 is evidence that they are fundamental to our ability to adapt, in real-time, to new people, new  
183 situations and new tasks. Conversation analysts have described a set of systematic turn-based *repair*  
184 processes that structure how people identify and respond to misunderstandings (Schegloff et al.,  
185 1977a; Schegloff, 1992a, 1997). Experimental evidence shows these repair processes have a critical  
186 role in building up shared understanding and shared languages on the fly (Healey et al., 2018b;  
187 Healey, 2008, 1997).

188 The Running Repairs Hypothesis characterises human communication as a fundamentally error-  
189 prone, effortful, active, collaborative process but also highlights how these processes are structured  
190 and how they make human communication flexible and adaptable to new people and new situations.  
191 This can liberate human-robot interaction from the fantasy of perfect competence (Park et al., 2021).  
192 Instead, robots could, in principle, take advantage of the resources of interaction by engaging in  
193 repairs. This requires developing the ability to recognise critical verbal and non-verbal signals of  
194 misunderstanding and the use of incremental online learning processes that build on the sequential  
195 structure of interaction to make real-time revisions to language models (see e.g. Howes and Eshghi  
196 2021; Purver et al. 2011).

### 197 3.1.2 Repair, Recruitment, and (virtual) Agency in a Smart Homecare Setting

198 Albert argued that moments of trouble and failure can provide researchers with ideal empirical  
199 material for observing the structure of the participation frameworks we use to get things done in  
200 everyday life (Goodwin, 2007; Albert and Ruiters, 2018). His presentation used multimodal video  
201 analysis to show how a disabled man and his (human) carer leveraged troubles and failures in their  
202 interactions with an Amazon Echo with voice-controlled lights, plugs, and other devices to co-design  
203 an effective smart homecare participation framework.

204 Instances in this case study highlighted how the human carer used troubles and failures to prioritise  
205 the independent role and agency of the disabled person within a joint activity. For example, the  
206 carer would stop and wait for the disabled person to resolve the trouble in their interactions with the  
207 virtual agent and complete their task even when it would have been faster for the carer to complete  
208 the disabled person's task manually. In other examples, trouble in the interactions between the carer  
209 and the virtual assistant provided an opportunity for the disabled person to intervene and assist  
210 the carer by correcting and completing their vocal instruction to the device. The disabled person  
211 was also able to tacitly 'recruit' (Kendrick and Drew, 2016) assistance from the human carer by  
212 repeatedly re-doing failed commands to the virtual assistant within earshot of the carer, soliciting  
213 support without having to ask for help directly.

214 These episodes show how people can harness trouble and failures in interaction with a virtual  
215 assistant to enable subtle shifts of agency and task-ownership between human participants. This  
216 kind of hybrid smart homecare setting can support and extend the independence of a disabled  
217 person within an interdependent, collaborative participation framework (Bennett et al., 2018). More  
218 broadly, the communicative utility of trouble and failure in interactions with machines highlights the  
219 shortcomings of our idealized—often ableist—models of the 'standard' user, and medicalized models  
220 of assistive technology (Goodwin, 2004; Albert and Hamann, 2021).

### 221 3.1.3 Building Common Ground in Human-robot Interaction

222 Skantze highlighted two aspects of miscommunication and error handling in human-machine  
223 interaction. First, he discussed how language is ultimately used as part of a joint activity.  
224 For communication to be meaningful and successful, the interlocutors need to have a mutual  
225 understanding of this activity, and of their common ground (Clark, 1996). From this perspective,  
226 language processing is not a bottom-up process, where we first figure out what is being said before  
227 interpreting and putting it in context. Rather, we use the joint activity to steer the interpretation  
228 process and possibly ignore irrelevant signals. Skantze exemplified this with an early experiment,  
229 where a noisy channel (including a speech recognizer) was used in a human-human communication  
230 task, where one person had to guide another person on a virtual campus (Skantze, 2005). Although  
231 much of what was said did not get through (due to the error prone speech recognition), the humans



232 very seldom said things like “sorry, I didn’t understand”, which are frequent responses in human-  
233 machine interactions. Instead, they relied on the joint activity to ask task-related questions that  
234 contributed to task progression. Another implication of this view on communication is that the  
235 idea of “open-domain dialogue”, where there is no clear joint activity, is not meaningful to pursue  
236 (Skantze and Doğruöz, 2023).

237 The second aspect that was discussed was the need to incorporate user feedback when the system  
238 is speaking, and use that feedback to model what can be regarded as common ground between the  
239 user and the system. Skantze exemplified this issue with a research project at KTH (Axelsson and  
240 Skantze, 2023), where an adaptive robot presenter is being developed (in the current demonstrator  
241 it is talking about classic works of art in front of a human listener). The robot presenter uses a  
242 knowledge graph to model the knowledge it is about to present, and then uses that same graph to  
243 keep track of the “grounding status” of the different pieces of information (Axelsson and Skantze,  
244 2020). Multimodal feedback from the user (e.g., gaze, facial expressions, nods and backchannels)  
245 are interpreted as negative or positive, and the graph is updated accordingly, so that the presentation  
246 can be adapted to the user’s level of knowledge and understanding (Axelsson and Skantze, 2022).

#### 247 3.1.4 Addressing the Challenges of Open-Domain Conversational AI Systems

248 Papaioannou’s presentation showed how designing conversational AI systems able to engage in  
249 open-domain conversation is extremely challenging and a frontier of current research. Such systems  
250 are required to have extensive awareness of the dialogue context and world knowledge, the user  
251 intents and interests, requiring more complicated language understanding, dialogue management,  
252 and state and topic tracking mechanisms compared to traditional task-oriented dialogue systems.

253 In particular, some of these challenges include: (a) keeping the user engaged and interested over  
254 long conversations; (b) interpretation and generation of complex context-dependency phenomena  
255 such as ellipsis and anaphora; (c) mid-utterance disfluencies, false starts, and self-corrections  
256 which are ever-present in spoken conversation (Schegloff et al., 1977b; Shriberg, 1994) (d) various  
257 miscommunication and repair phenomena such as Clarification Requests (Purver, 2004) and Third  
258 Position Repair (Schegloff, 1992b) whereby either the user or system does not understand the other  
259 sufficiently or misunderstands, and later repairs the misunderstanding. (b-d) are all crucial to robust  
260 Natural Language Understanding in dialogue.

261 A modular conversational AI system, (called *Alana*), tackling some of the aforementioned  
262 challenges (i.e. user engagement over long conversations, ellipsis and anaphora resolution, and  
263 clarification requests) was developed between 2017-2019 (Papaioannou et al., 2017; Curry et al.,  
264 2018) and deployed to thousands of users in the United States as part of the Amazon Alexa Challenge  
265 (Ram et al., 2018). The Alana system was also evaluated in a multimodal environment and was used  
266 as the overall user conversational interaction module in a multi-task and social entertainment robotic

267 system as part of the MuMMER project (Foster et al., 2019). The integrated system was deployed in  
268 a shopping mall in Finland and was able to help the user with specific tasks around the mall (e.g.  
269 finding a particular shop or where they could buy a certain product, finding the nearest accessible  
270 toilet, or asking general questions about the mall) while at the same time engaging in social dialogue  
271 and being entertaining.

272 The output of that research was fed to the implementation of the ‘Conversational NLU’ pipeline by  
273 Alana AI, a modular neuro-symbolic approach further enhancing the language understanding of the  
274 system. The Conversational NLU module is able to detect and tag a number of linguistic phenomena  
275 (e.g. disfluencies, end-of-turn, anaphora, ellipsis, pronoun resolution, etc) as well as detect and  
276 repair misunderstandings or lack of sufficient understanding, such as self-repairs, third-position  
277 corrections, and clarifications. The system is currently being evaluated by blind and partially sighted  
278 testers in the context of multi-modal dialogue allowing the users to find mislocated objects in their  
279 environment via a mobile application.

## 280 **3.2 Lightning Talks**

281 The following section contains short summaries of the lightning talks of both the virtual and the  
282 face-to-face part of the workshop. From the presentations, three themes were identified: *Description*  
283 *and Analysis of Failures and Troubles* (Sect. 3.2.1) grouping presentations that have a descriptive  
284 or analytical focus; *Technical Aspects of Conversational Failure* (Sect. 3.2.2) for presentations  
285 that have a more technical focus; and *Adjacent Topics in Speech Interfaces* (Sect. 3.2.3), grouping  
286 presentations on topics that, while not focusing strictly on conversational failures, covering other  
287 forms of errors and issues that fall into the wider topic of speech-centric human-machine interactions.  
288 Note that many of the talks falling into the second, technical category still contain a substantial  
289 element of analysis that enabled or inspired the technical solutions described therein.

### 290 **3.2.1 Description and Analysis of Failures and Troubles**

291 The following ten of the contributions took a more analytical approach to the failure they reported  
292 in their lightning talks. They describe possible reasons or implications of the failure they present.

#### 293 **3.2.1.1 Laundrobot: Learning from Human-Human Collaboration**

294 Barnard and Berumen presented their work on *Laundrobot*, a human acting as a collaborative robot  
295 designed to assist people in sorting clothing into baskets. The study focused on participants’ ability  
296 to collaborate through verbal instructions and body movements with a robot that was sometimes  
297 erroneous when completing the task. The team analysed social signals, including speech and gestures,  
298 and presented three cases demonstrating human-human collaboration when things do not go as  
299 expected. In one of the cases, a participant gave clear instructions to an erroneous Laundrobot, which

300 led to frustration on the participant's part, with statements such as "Okay, I'm doing this wrong".  
 301 The presenters described how the participant appeared to take responsibility for the errors made by  
 302 the robot. They examined the use of language and expression of intent in different instances for  
 303 pieces of clothing that were either correctly or incorrectly identified by Laundrobot. During this  
 304 analysis, Barnard, Berumen, and colleagues came across an interesting case regarding the use of the  
 305 word "right", which was frequently used in both erroneous and non-erroneous instances. The group  
 306 explored how that word had different meanings depending on the success or failure of Laundrobot.  
 307 For instance, for one participant (P119), the word had a single meaning of indicating a direction in  
 308 erroneous instances, whereas, on other occasions, it had alternative purposes. It was sometimes used  
 309 to refer to directions and, at other times, used for confirmation, immediacy ("right in front of you"),  
 310 or purpose ("Right, OK").

### 311 **3.2.1.2 Sequential Structure as a Matter of Design and Analysis of Trouble**

312 As part of the *Peppermint project*<sup>3</sup> corpus, Tisserand presented a transcript fragment, reproduced  
 313 below. They designed a Pepper robot as an autonomous reception desk agent that would answer  
 314 basic requests asked by library users. They captured *naturally-occurring interactions*: the robot was  
 315 placed in the library, and users were free to interact and leave whenever they wanted.

316	01 Hum: where can I find books of maths?	Sequence A - Part 1
317	02 Rob: ((provides the direction for books of maths))	Sequence A - Part 2
318	03 Rob: is it clear to you?	Sequence B - Part 1
319	04 Hum: yes thanks	Seq B-2 && Seq A-3
320	05 Rob: okay, I will repeat ((repeats turn line 2))	Sequence C - Part 1

321 The failure here is the fact that the robot recognized "no thanks" instead of two separate actions:  
 322 "yes" + "thanks" (l.4); the robot thus repeats the answer to the user's question. Reflecting on this  
 323 WTF moment, Tisserand highlighted how this failure occurred due to decisions made during the  
 324 scenario design phase. Firstly, poor speech recognition differentiation between the words "yes" and  
 325 "no" had led the scenario design team to add "no thanks" to a word list provided for recognising  
 326 an *offer rejection*: (a *dispreferred turn design* for this type of action (Schegloff, 2007, Chap.5)) in  
 327 another scenario in which the robot makes an offer. Secondly, because the state machine was based  
 328 on isolated so-called "contexts", it was designed only to make one decision when processing a spate  
 329 of talk. Here, therefore, the clarification check turn in line 3 was treated as independent from the  
 330 question response in line 2. Because the speech recognition system struggled to differentiate "yes"  
 331 and "no", and was using the word list that labelled "no thanks" as a case of *offer rejection*, here it  
 332 erroneously recognized "yes thanks" in line 4 as a negation (a *clarification denial*), and proceeded  
 333 to repeat the turn.

<sup>3</sup> <https://peppermint.projet.liris.cnrs.fr/>

334 What should have happened is that when the robot asks the user to confirm (1.3), it should recognize  
335 that this sequence is embedded in the previous question/answer sequence (1.1-2). In this case, the  
336 human's "yes" (1.3) is a response to the just-prior confirmation request while the "thanks" responds  
337 (in the first structurally provided sequential slot) to the Robot's answer as a 'sequence closing third'  
338 (1.3). This is why the team is now *sequentially* annotating training datasets to show what utterances  
339 correspond not only to questions and answers, but also the cement in-between: how the user might  
340 delay, suspend, abandon, renew or insert actions (e.g. repair). Here interaction is seen as a temporally  
341 continuous and incremental process and not a purely logical and serial one. In other words, context  
342 is seen as an organized resource more than an adaptability constraint.

### 343 **3.2.1.3 Design a Robot's Spoken Behaviours Based on How Interaction Works**

344 Huang pointed out that spoken interaction is complicated. It is grounded in the social need to  
345 cooperate (Tomasello, 2009; Holtgraves, 2013) and requires interlocutors to coordinate and build  
346 up common ground on a moment-by-moment basis (Krauss and Fussell, 1990, p.112)(Holtgraves,  
347 2013).

348 Speech is only one tool in a larger picture. Some errors are caused by failures in natural language  
349 understanding (NLU) as illustrated in the following sequence:

350 01 User: Let's talk about me.  
351 02 Robot: What do you want to know about 'me'?

352 Other issues, however, could be caused by a lack of understanding of common ground. For example,  
353 when a naive user asked, "Where to find my Mr Right", the system provided a place named "Mr  
354 & Mrs Right" and told the user it was far away. This reply contains several layers of failure: (1)  
355 the robot fails to capture the potential semantic inference of the expression *Mr Right*; (2) it fails  
356 to consider the social norm that Mr Right belongs typically to one person only; and (3) it makes  
357 a subjective judgement about distance. One may argue that this error would not happen if the  
358 user knew a question-answer robot could not chat casually. However, the issue is whether a clear  
359 boundary of a social robot's capability is set in the system or communicated to the user during the  
360 interaction. It is difficult to tell why speech interfaces may fail and how to work around the limits  
361 without understanding what makes interaction work and how speech assists in the process.

362 Also, spoken interaction requires interlocutors, including robots, to adjust their behaviours based  
363 on the verbal and non-verbal feedback provided by others. A social robot that does not react  
364 appropriately could be deemed improperly functional, as illustrated in the following sequence. In  
365 the scenario, the robot failed to generate satisfactory answers several times in an open conversation;  
366 the user felt frustrated.

367 User: You are generating GPT rubbish.

368 Robot: (No response, carries on)

### 369 **3.2.1.4 Hey Siri ... You Don't Know How to Interact, huh?**

370 The WTF moment Wiltchko presented concerned the use of *huh* in interaction with Siri, Apple's  
371 voice assistant.

372 User: Hey Siri, send an e-mail.

373 Siri: To whom shall I send it?

374 User: huh?

375 Siri: I couldn't find huh in your contacts. To whom shall I send it?

376 It is evident from the example that Siri cannot understand *huh*. This is true for *huh* used as an  
377 other-initiated repair strategy as in the example above, but it is also true for its use as a sentence-final  
378 tag. This is a significant failure as in human-human interaction the use of *huh* is ubiquitous. In fact,  
379 *huh* as a repair strategy has been shown to be available across a number of unrelated languages  
380 (Dingemanse et al., 2013). Wiltchko speculates that successful language use in machines is restricted  
381 to propositional language (i.e., language used to convey content) whereas severe problems arise in  
382 the domain of interactional language (i.e., language used to regulate common ground building as  
383 well as the conversational interaction itself). The question that arises, however, is whether human  
384 users feel the need to use interactional language with machines. After all, this aspect of language  
385 presupposes interaction with another mind for the purpose of common ground construction and it  
386 is not immediately clear whether humans treat machines as having a mind with which to share a  
387 common ground.

### 388 **3.2.1.5 Utilising Explanations to Mitigate Robot Failures**

389 Kontogiorgos presented current work on failure detection (Kontogiorgos et al., 2020a, 2021)  
390 and how robot failures can be used as an opportunity to examine robot explainable behaviours.  
391 Typical human-robot interactions suffer from real-world and large-scale experimentation and tend to  
392 ignore the 'imperfectness' of the everyday user (Kontogiorgos et al., 2020b). Robot explanations  
393 can be used to approach and mitigate robot failures by expressing robot legibility and incapability  
394 (Kwon et al., 2018), and within the perspective of common-ground. The presenter discussed  
395 how failures display opportunities for robots to convey explainable behaviours in interactive  
396 conversational robots according to the view that miscommunication is a common phenomenon  
397 in human-human conversation and that failures should be viewed as being an inherent part of  
398 human-robot communication. Explanations, in this view, are not only justifications for robot actions,  
399 but also embodied demonstrations of mitigating failures by acting through multi-modal behaviours.

### 400 **3.2.1.6 Challenging Environments for Debugging Voice Interactions**

401 Porcheron presented the challenge of how we expect users to understand and debug issues with  
402 ‘eyes-free voice interactions’, and of parallelism to the prospects of voice-based robots. A recurrent  
403 promise of voice-based technologies is their simplicity: we issue a command to a computer and it can  
404 respond accordingly. Of course, not all technology use goes as planned and sometimes errors occur.  
405 With graphical user interfaces (GUIs), we have a plethora of well-tested heuristics (e.g., Nielsen  
406 (1995)), especially for dealing with ‘errors’ where users need ‘fix’ something. However, with voice,  
407 in situations where people encounter something going wrong, they have to carry out work to figure  
408 out how to resolve the issue (Porcheron et al., 2018; Fischer et al., 2019). One specific example  
409 is responses which do not reveal specifics, such as “I had an issue responding to that request”.  
410 Users are given little purchase with which to debug this issue, and attempt to resolve this. This user  
411 challenge is exacerbated in the new settings where voice technologies are appearing: in our cars,  
412 on our bikes, and anywhere we take our smartwatch—in these settings, there is often little time  
413 to read and respond to a text, little audible information to go on, and plenty of distraction for the  
414 user. Porcheron suggested that if we want to consider voice as a modality for controlling robots, we  
415 first need to think through how we help users understand and recover from ‘errors’ in these sorts of  
416 environments first.

### 417 **3.2.1.7 The Robot [does not] Hear Me Correctly**

418 Hernández García presented examples of failures of user interactions with a receptionist robot.  
419 Gunson et al. (2022) developed a social conversational robot that can act as a receptionist in a  
420 hospital waiting room for tasks such as greeting visitors, answering common questions, providing  
421 directions and/or scheduling information, etc. During the evaluation with users (student participants  
422 in a lab setting), a majority of failed interactions that prevented user task completion were the result  
423 of problems with the ASR. These ASR failures predominately took the form of an ASR result that  
424 contained only part of the user’s utterance, misrecognition of critical words, or a combination of  
425 the two (to the extent that the user’s original intention could not be discerned), leading the dialogue  
426 management into a non-coherent path and frustrating the user.

427 USER: can you help me to find the exit please | ASR: please  
428 ROBOT: Anyway it’s a big building we’re in, isn’t it.  
429 Did you know you can ask me for help if you ever get lost here?  
430 USER: can you direct me to the exit please | ASR: can you direct me  
431 to Netflix please  
432 ROBOT: Such a lot of waiting around isn’t there.  
433 USER: thank you [user left]

434 Participants rated the conversation negatively when these errors occurred, as they experienced  
435 difficulties in making themselves understood. The user evaluations reported by Gunson et al. (2022)  
436 highlighted that users did not feel it was *natural* or that it *flowed* in the way they expected. Participants  
437 did not believe that “*the robot heard me correctly most of the time*” or that “*the robot recognised the*  
438 *words I said most of the time*” nor “*felt confident the robot understood the meaning of my words*”.

439 Conversational troubles may start at a *speech recognition* level, but these failures are propagated  
440 throughout the whole *speech interface* pipeline, compounding to create WTF moments and leading  
441 to poor performance, increasing user frustration, and loss of trust, etc.

### 442 **3.2.1.8 Hello, It’s Nice to “Meat” You**

443 Nettet shared examples of WTF moments encountered while interacting with Norwegian chatbots  
444 through written text. The first failure presented was users’ committing spelling mistakes interacting  
445 with a virtual agent through chat. This caused the agent to misunderstand the overall context of the  
446 conversation. A good example of this is misspelling meet with meat, and the chatbot then replying  
447 with a response about sausages.

448 The second part entailed a user failure that is specifically for multilingual users. In some non-native  
449 English-speaking countries, such as Norway, technical terms and newer words are often commonly  
450 said in English. This potentially leads users to interact with agents in two languages within the same  
451 sentence/conversation. This can lead to the agent struggling to interpret the terms in the second  
452 language, and assuming that they mean something else in the original interaction language. These  
453 are some examples of how uncertain user output can result in failures from the robot.

### 454 **3.2.1.9 Speech Misrecognition: A Potential Problem for Collaborative Interaction in** 455 **Table-grape Vineyards**

456 Kaszuba presented troubles and failures encountered while designing a spoken human-robot  
457 interaction system for the *CANOPIES project*<sup>4</sup>. This project aims to develop a collaborative paradigm  
458 for human workers and multi-robot teams in precision agriculture, specifically in table-grape  
459 vineyards. When comparing some already existing speech recognition modules (both online and  
460 offline), the presenter identified communication issues associated with the understanding and  
461 interpretation of specific words of the vineyard scenario, such as “grape”, “bunch”, and “branch”.  
462 Most of the tested applications could not clearly interpret such terms, leading the user to repeat the  
463 same sentence/word multiple times.

464 Hence, the most significant source of failure in speech interfaces that Kaszuba has described is  
465 *speech misrecognition*. Such an issue is particularly relevant, since the quality and effectiveness of

---

<sup>4</sup> <https://www.canopies-project.eu/>

466 the interaction strictly depend on the percentage of words correctly understood and interpreted. For  
467 this reason, the choice of the application scenario has a crucial role in the spoken interaction, and  
468 preliminary analysis should be taken into consideration when developing such systems, as the type  
469 and position of the acquisition device, the ambient noise and the ASR module to adopt. Nevertheless,  
470 misrecognition and uncertainty are unavoidable when the developed application requires people  
471 to interact in outdoor environments and communicate in a language that is not the users' native  
472 language.

473 Hence, some relevant considerations concerning ASR modules should be taken into account in  
474 order to implement a robust system that, eventually, can also be exploited in different application  
475 scenarios. The percentage of uncertainty, the number of misrecognized words and the environmental  
476 noise that can negatively affect communication are some fundamental issues that must be addressed  
477 and minimized.

### 478 **3.2.1.10 Leveraging Multimodal Signals in Human Motion Data During Miscommunication** 479 **Instances**

480 Approaching from a natural dialogue standpoint and inspired by the Running Repairs Hypothesis  
481 Healey et al. (2018b), Özkan shared a presentation on why and how we should take advantage of  
482 WTF-moments or miscommunications to regulate shared understanding between humans and speech  
483 interfaces. Rather than avoiding these moments (which is impossible), if speech interfaces were to  
484 identify them and show appropriate behaviour, it could result in more natural, dynamic and effective  
485 communication.

486 Detecting miscommunications from the audio signal can only can be costly in terms of  
487 computational load or prone to error due to noise in most environments. Fortunately, repair  
488 phenomena manifest themselves in non-verbal signals as well Healey et al. (2015); Howes et al.  
489 (2016). Findings regarding speaker motion during speech disfluencies (self-initiated self-repairs)  
490 have shown that there are significant patterns in the vicinity of these moments Özkan et al. (2021,  
491 2023); Ozkan et al. (2022). Specifically, the speakers have higher hand and head positions and  
492 velocities near disfluencies. This could be treated as a clear indicator for artificial interfaces to  
493 identify troubles of speaking in their human partner. For example, to the user input "*Could you*  
494 *check the flights to Paris -uh, I mean- Berlin?*", the interface, instead of disregarding the uncertain  
495 utterance, could offer repair options more actively by returning "*Do you mean Paris or Berlin?*" in  
496 a collaborative manner.

497 Though not in the context of disfluencies, a common example of not allowing repair (in this case  
498 other initiated other repair) occurs when the user needs to correct the output of an interface or  
499 simply demand another response to a given input. As a WTF moment in the repair context, Özkan  
500 demonstrated a frequent problem in their interaction with Amazon Alexa. When asked to play a



501 certain song, Alexa would play another song with the same or similar name. The error is not due to  
502 speech recognition, because Alexa understands the name of the song very well. However, it maps  
503 the name to a different song that the user does not want to hear. No matter how many times the  
504 user tries the same song name input, even with the artist name, Alexa would still pick the one that  
505 is the ‘first’ result of its search. If the conversational repair was embedded in the design, a simple  
506 solution to this problem could have been “*Alexa, not that one, can you try another song with the*  
507 *same name?*”, but Alexa does not respond to such requests.

### 508 3.2.2 Technical Aspects of Conversational Failure

509 The following five of the contributions describe technical aspects of failures. Presentations in this  
510 section either discuss the technical causes of failures, point out technological attempts to recognize  
511 when conversational trouble occurs, or summarize approaches on handling troubles on part of the  
512 robot.

#### 513 3.2.2.1 *Chefbot: Reframing Failure as a Dialogue Goal Change*

514 Gkatzia presented their work on *Chefbot*, a cross-platform dialogue system that aims to help users  
515 prepare recipes (Strathearn and Gkatzia, 2021a). The task moves away from classic instruction  
516 giving and incorporates question-answering for clarification requests, and commonsense abilities,  
517 such as swapping ingredients and requesting information on how to use or locate specific utensils  
518 (Strathearn and Gkatzia, 2021b). This results in altering the goal of the communication from cooking  
519 a recipe to requesting information on how to use a tool, and then returning to the main goal. It  
520 was quickly observed that changing the dialogue goal from completing the recipe to providing  
521 information about relevant tasks resulted in failure of task completion. This issue was subsequently  
522 addressed by *reframing* failure as a temporary dialogue goal change, which allowed the users to  
523 engage in question answering that was not grounded to the recipe document, and then forcing the  
524 system to resume the original goal.

#### 525 3.2.2.2 *Failure in Speech Interfacing with Local Dialect in a Noisy Environment*

526 Liza (Farhana) presented their ongoing work in capturing the linguistic variation of speech  
527 interfaces in real-world scenarios. Specifically, local dialects may impose challenges when modelling  
528 a speech interface using an artificial intelligence (deep learning) language modelling system. Deep  
529 learning speech interfaces rely on language modelling which is trained on large datasets. A large  
530 dataset can capture some linguistic variations; however, dialect-level variation is difficult to capture  
531 as a large enough dataset is unavailable. Moreover, very large models require high-performance  
532 computation resources (e.g., GPU) and take a long time to respond, which imposes further constraints  
533 in terms of deploying such systems in real scenarios. Large data-driven solutions also cannot easily  
534 deal with noise as it is impractical to give access to enough real-world data from noisy environments.

535 Overall, state-of-the-art AI models are still not deployable in scenarios with dialect variation and  
536 noisy environments. Alharbi et al. (2021) identified several hurdles in training end-to-end Automatic  
537 Speech Recognition (ASR) models. Additionally, the conditional interdependence between the  
538 acoustic encoder and the language model was emphasized by (Xu et al., 2020). Consequently, while  
539 augmenting the standard text training data can enhance the efficacy of general-purpose language  
540 models, the limited availability of corresponding acoustic data poses challenges in training end-to-  
541 end ASR systems. Moreover, when addressing dialect modeling (Hirayama et al., 2015), the scarcity  
542 of training data exacerbates the difficulties in integrating speech interfacing and language modeling  
543 (Liza, 2019) within the ASR framework.

### 544 **3.2.2.3 The ‘W’ in WTF Moments can also be ‘When’: The Importance of Timing and** 545 **Fluidity**

546 Hough presented WTF moments driven more by inappropriate timing of responses to user  
547 utterances, rather than by content misunderstandings. Improving the first-time accuracy of Spoken  
548 Language Understanding (SLU) remains a priority for HRI, particularly given errors in speech  
549 recognition, computer vision and natural language understanding remain pervasive in real-world  
550 systems. However, building systems capable of tolerating errors whilst maintaining *interactive*  
551 *fluidity* is an equally important challenge. In human-human situated interactions where an instructee  
552 responds to a spoken instruction like “put the remote control on the table” and a follow-up repair  
553 like “no, the left-hand table” when the speaker realizes the instructee has made a mistake, there is  
554 no delay in reacting to the initial instruction, and adaptation to the correction is instant (Heldner  
555 and Edlund, 2010; Hough et al., 2015), in stark contrast to state-of-the-art robots with speech  
556 interfaces. Increasing interactive fluidity is vital to give robots with speech understanding more  
557 seamless, human-like transitions from processing speech to taking physical action without delay,  
558 permitting appropriate overlap between the two, and the ability to repair actions in real-time. Rather  
559 than waiting for components to be perfected, preliminary experiments with a pick-and-place robot  
560 show users can be tolerant of errors if fluidity is kept high, including appropriate repair mechanisms  
561 (Hough and Schlangen, 2016).

### 562 **3.2.2.4 Laughter in WTF Moments**

563 Maraev presented a hypothesis that laughter can be treated as an indicator of a WTF moment.  
564 Laughter can occur in such moments as a) speech recognition failures disclosed to a user via explicit  
565 grounding feedback, b) awkwardness due to retrieval difficulties, c) resulting system apologies and  
566 down players (e.g., “don’t worry”). Along with examples from task-oriented role-played dialogues,  
567 Maraev discussed the following constructed example, where laughter communicates a negative  
568 feedback to the system’s clarification of speech recognition result:

569 Ustr> I would like to order a vegan bean burger.  
570 Sys> I understood you'd like to order a vegan beef burger. Is that correct?  
571 Ustr> HAHAHA

572 Maraev et al. (2021) focused on non-humorous laughs in task-oriented spoken dialogue systems.  
573 The paper shows how certain types of laughter can be processed within the dialogue manager and  
574 natural language generator, namely: laughter as negative feedback, laughter as a negative answer to  
575 a polar question and laughter as a signal accompanying system feedback.

### 576 **3.2.2.5 To Err is Robot**

577 Giuliani presented findings from six years of research on erroneous human-robot interactions.  
578 The team of researchers led by Giuliani has shown that participants in human-robot interaction  
579 studies show unique patterns of social signals when they experience an erroneous situation with  
580 a robot (Mirnig et al., 2015). The team annotated two large video corpora of 201 videos showing  
581 578 erroneous situations and 1200 videos showing 600 erroneous situations, respectively (Giuliani  
582 et al., 2015; Cahya et al., 2019). They found that there are two types of errors that do occur in  
583 human-robot interaction. Social norm violations are situations in which the robot does not adhere  
584 to the underlying social script of the interaction. Technical failures are caused by the technical  
585 shortcomings of the robot. The results of the video analysis show that the study participants use  
586 many head movements and very few gestures but they often smile when in an error situation with  
587 the robot. Another result is that the participants sometimes stop moving at the beginning of error  
588 situations. The team was also able to show in a user study for which a robot was purposefully  
589 programmed with faulty behaviour that participants liked the faulty robot significantly better than  
590 the robot that interacted flawlessly (Mirnig et al., 2017). Finally, the team trained a statistical model  
591 for the automatic detection of erroneous situations using machine learning (Trung et al., 2017). The  
592 results of this work demonstrate that automatic detection of an error situation works well when the  
593 robot has seen the human before.

### 594 **3.2.3 Adjacent Topics in Speech Interfaces**

595 The two contributions under this theme do not discuss conversational failures directly but address  
596 the related topics of explanatory AI and privacy of speech interfaces.

#### 597 **3.2.3.1 What is a 'Good' Explanation?**

598 Kapetanios presented some thoughts around the long-standing research question of *what is a*  
599 *good explanation* in the context of the current buzz around the topics of explainable AI (XAI)  
600 and interpretable Machine Learning (IML). Using Amazon's Alexa and Google's Digital Assistant  
601 to generate explanations for answers being given to questions being asked of these systems, he

602 demonstrated that both systems, at the technological forefront of voice-based HCI approaches to  
603 answering specific questions, fail to generate convincing explanations. Convincing explanations  
604 should fit the facts, be relevant, tailored to the recipient, and typically do more than merely describe  
605 a situation (Dowden, 2019, chap. 14). It is frequently the latter where digital assistants have been  
606 observed to struggle. Hence, when describing the results of running several thousand queries through  
607 the most common digital assistants, provides the following example (Enge, 2019):  
608 Siri, when being asked the question “Who is the voice of Darth Vader?”, instead of providing  
609 the name of the (voice) actor, returns a list of movies featuring Darth Vader. While this answer  
610 is topically relevant, it certainly is not a proper answer to the question. The same problem of  
611 explanation persists with ChatGTP-3/4, despite its fluency in generating precise answers to specific  
612 questions in natural language.

### 613 **3.2.3.2 Privacy and Security Issues with Voice Interfaces**

614 Williams presented privacy and security issues and how these are often underestimated, overlooked,  
615 or unknown to users who interact with voice interfaces. What many voice interface users are unaware  
616 of is that only three to five seconds of speech are required to create a *voiceprint* of a person’s real  
617 voice as they are speaking (Luong and Yamagishi, 2020). One of the risks that follows is that  
618 voiceprints can be re-used in other voice applications to impersonate or create voice deepfakes  
619 (Williams et al., 2021b,a). In the UK and many other countries, this poses a particular security risk  
620 as voice-authentication is commonly used for telephone banking and call centres. In addition, some  
621 people may be alarmed when a voice interface reveals private information by “speaking out loud”  
622 sensitive addresses, birth dates, account numbers, or medical conditions. Anyone in the nearby  
623 vicinity may overhear this sensitive information and technology users have no ability to control what  
624 kinds of information a voice interface may say aloud (Williams et al., 2022).

### 625 **3.2.4 Summary of Lightning Talks**

626 Through their lightning talks, our participants contributed to an initial gathering of different  
627 troubles and failures in conversational interactions between humans and robots. Thanks to the  
628 description of their memorable failures and their analysis, we could identify the themes of *analysis*,  
629 *technical aspects* and *adjacent topics*, which all impact the success (or failure) of a conversation.

### 630 **3.3 Summary of World Café Session**

631 During the World Café session, four working groups were created based on recurring themes  
632 from the lightning talks, participants’ answers as to what they perceived as the most pressing issue  
633 or the biggest source of failure for speech interfaces, as well as the aim to define the sought after  
634 benchmark scenario. Through the initial submissions of the participants, their lightning talks and the  
635 keynotes, three main macro-categories have emerged: i) miscommunication, ranging from speech

636 recognition failures to more semantic and conversation-dependent failures; ii) interaction problems,  
637 encompassing all those failures that are due to users' expectations and behaviours; iii) context  
638 understanding, linked to the fact that interaction is shaped by context and that context changes fast,  
639 calling for a need to find more robust ways to establish common ground. While these three themes  
640 are highly interdependent and could culminate in the sought after benchmark scenario (the fourth  
641 working group), each of them presents peculiarities that we considered worth discussing in detail.

### 642 3.3.1 Handling Miscommunication

643 The discussion focused on the need to acknowledge and embrace the concept of miscommunication.  
644 One of the open challenges identified by this group was to equip robots with the ability to learn  
645 from various forms of miscommunication and to actively use them as an opportunity to establish  
646 common ground between users and robots. When communicating with a robot, the human user  
647 usually has a goal in mind. The robot could exploit miscommunication to understand this goal  
648 better by asking for clarifications at the right moments and updating the common ground. The  
649 discussion also acknowledged that miscommunication is only the starting point. Two distinct new  
650 challenges and opportunities arise when working on resolving miscommunication: 1) how to explain  
651 the miscommunication, and 2) how to move the conversation forward. Both problems are highly  
652 context-dependent and related to the severity and type of miscommunication. Moreover, being  
653 able to repair a breakdown in conversation may also depend on being able to establish appropriate  
654 user expectations in the first place by giving an accurate account of what the robot is really able  
655 to accomplish. The final discussion point from this group centered on the possibility of enriching  
656 the multimodal and non-verbal component of conversations to help the robot perceive when a  
657 miscommunication has happened by detecting and responding to, for example, long pauses or  
658 changes in specific types of facial expressions.

### 659 3.3.2 Interaction Problems

660 Interaction problems do not only encompass challenges that are specific to the technology used,  
661 like issues with automatic speech recognition or the presence of long delays when trying to engage  
662 in a "natural" conversation. They are related to perceived failures that longitudinally include all the  
663 technical problems identified by the other themes and relate to how the interaction with the human  
664 user is managed. In this context, human users play an essential role and the participants of this  
665 group emphasized the necessity of creating expectations that allow users to build an adequate mental  
666 model of the technology they are interacting with. In Washburn et al. (2020a), authors examine how  
667 expectations for robot functionality affected participants' perceptions of the reliability and trust of a  
668 robot that makes errors. The hope is that this would lead to an increased willingness and capacity  
669 to work with the failures that inevitably occur in conversational interactions. Anthropomorphism  
670 was identified as one of the possible causes for the creation of wrong expectations: the way robots

671 both look and speak risks tricking users into thinking that robots have human-like abilities and are  
672 able to follow social norms. Once this belief is abandoned, users could then form an appropriate  
673 expectation of the artificial agents, and the severity of the failures would decrease. Setting the right  
674 expectations will also enable users to understand when a failure is a technological error in execution  
675 or when it is a design problem: humans are unpredictable, and some of the problems that arise in the  
676 interactions are due to users' behaviours that were not embedded in the design of robot's behaviours.  
677 A related aspect that was considered important by this group is the transparency of the interaction:  
678 the rationale behind the failures should be explained and made clear to the users to enable mutual  
679 understanding of the situation and prompt recovery. This could, in fact, be initiated by the users  
680 themselves. Another need, identified as a possible way to establish better conversational interactions,  
681 is the missing link of personalisation. The more the agents are able to adapt to the context and the  
682 users they are interacting with, the more they will be accepted, as acceptance plays a fundamental  
683 role in failure management. A general consensus converged regarding the fact that we are not yet  
684 at the stage where we can develop all-purpose chatbots - or robots - and the general public should  
685 be made aware of this, too. Each deployment of conversational agents is context related and the  
686 conversation is mainly task-oriented, where a precise exchange of information needs to happen for a  
687 scenario to unfold.

### 688 3.3.3 Context Understanding

689 All four groups agreed that context understanding is crucial for reducing or entirely eliminating  
690 failures of interactive systems that use spoken language. We determined that capturing and modelling  
691 context is particularly challenging since it is an unbound and potentially all-encompassing problem.  
692 Moreover, all dialogue, and in fact, interaction as a whole, would be *shaped by* the context while at  
693 the same time *renewing* it. Likewise, the volatility of context, in particular, potentially rapid context  
694 switches, was also identified as challenging in human-robot conversation. Modelling the interaction  
695 partner(s) and evaluating their focus of attention was thereby discussed as one potential approach to  
696 reducing context search space.

697 A precise and consistent representation of the dialogue context was therefore identified as one of  
698 the most important problems that would rely on modelling not only the current situation but also any  
699 prior experiences of humans with whom the system is interacting. Such previous experience was seen  
700 to have significant effects on expectations about the interactive system that would potentially require  
701 calibration before or during system runtime to avoid misunderstandings as well as misaligned trust  
702 towards the system Hancock et al. (2011). However, even if we assume an optimal representation of  
703 context would be possible, the problem of prioritisation and weighting would still persist.

704 Another challenge discussed was the need for a multi-modal representation of the current situation  
705 comprised of nonverbal signals, irregular words, and interjections. Such a model would be required

706 for an appropriate formulation of common ground, whereby it remains unclear what exactly would  
707 be required to include. In that context, one group identified the benefits of a typology that could  
708 encompass an interaction situation in a multi-modal way, potentially extending work by Holthaus  
709 et al. (2023). The exact mapping between a signal or lexical index and their meanings is, however,  
710 still difficult to establish.

711 On the other hand, considering the dialogue context was unanimously regarded as beneficial to  
712 enrich human-robot conversations offering numerous opportunities to increase its functionality, even  
713 if it would not be possible to capture all context comprehensively. With a personalised model of  
714 interaction partners, for example, the spoken dialogue could be enhanced by taking into account  
715 personal interaction histories and preferences. Conversational agents could be improved for highly  
716 constrained settings and converge faster to relevant topics.

717 It is noteworthy to mention that enriching the capabilities of conversational agents with context  
718 information poses ethical challenges, e.g. in terms of privacy and data protection. This approach  
719 might thus introduce barriers in terms of user acceptance that need to be considered Lau et al. (2018).  
720 However, using context appropriately could also help to improve a system's transparency either by  
721 designing it with its intended context in mind or by utilising it during a conversation, for example,  
722 by providing additional interfaces to transport further information supporting the dialogue or by  
723 analysing context to reduce ambiguities and eliminate noise. The context was regarded to often play  
724 a vital role in providing the necessary semantic frame to determine the correct meaning of spoken  
725 language. Making use of domain and task knowledge was thereby identified as particularly helpful.

726 Moreover, intentionally misapplying context or analysing situations where context has previously  
727 misled a conversation, might be avenues to recognize and generate error patterns to help detect  
728 future troubles and failures in speech understanding.

#### 729 3.3.4 Benchmark Scenario(s)

730 On this discussion table, participants struggled to devise a single benchmark scenario that would  
731 elicit most, if not all, commonly occurring conversational failures. As a main reason for the difficulty  
732 of identifying such a prototypical scenario, the lack of a comprehensive taxonomy of conversational  
733 failures was determined.

734 An alternative suggestion to the proposed task of identifying one, failure-wise all encompassing,  
735 scenario was also made. Rather than seeking to specify a single scenario, it may be necessary  
736 to create test plans for each specific interaction task using chaos engineering, with some of the  
737 defining characteristics for a scenario being (1) the type(s) of users, (2) the domain of use (e.g.  
738 health-related, shopping mall information kiosk), (3) the concrete task of the robot, (4) the types  
739 of errors under investigation. Chaos engineering is typically used to introduce a certain level of  
740 resilience to large distributed systems (cf. Fomunyan (2020)). Using this technique, large online

741 retailers such as Amazon deliberately knock out some of their subsystems, or introduce other kinds  
742 of errors, to ensure that the overall service can still be provided despite the failure of one or more  
743 of these, typically redundant, components (cf. Siwach et al. (2022)). While both the envisioned  
744 benchmark scenario(s) and chaos engineering are meant to expose potential failures of human-made  
745 systems, the types of systems and types of failure differ substantially. While failures in technical  
746 distributed systems are unilateral, in the sense that the source of failure is typically attributed solely  
747 to the system rather than its user, attribution of blame in conversational failure is less unilateral. If a  
748 successful conversation is seen to be a joint achievement of at least two speakers, conversational  
749 failure is probably also best seen as a joint “achievement” of sorts. In other words, the *user* of a  
750 conversational robot is always also an interlocutor during the interaction. Hence, whatever approach  
751 we use to identify and correct conversational failures, the correct level of analysis is that of the dyad  
752 rather than of the robot alone.

753 Independent of the chaos engineering approach, another suggestion was that at least two benchmarks  
754 might be needed in order to distinguish between low-risk and high-risk conversations. Here, low-risk  
755 conversations would be the more casual conversations that one may have with a shop assistant whose  
756 failure would not carry any hefty consequences. High-risk conversations, on the other hand, would  
757 be those where the consequences of conversational failure might be grave - imagine conversational  
758 failure between an assistive robot and its human user that are engaged in some joint task of removing  
759 radioactive materials from a decommissioned nuclear site. If such a distinction should be made, the  
760 logical follow-up question would be how the boundary between low and high-risk scenarios should  
761 be determined. Finally, it should be mentioned that at least partial benchmarks such as *Paradise*  
762 exist for the evaluation of spoken dialogue systems Walker et al. (1997).

## 4 DISCUSSION

763 One significant result from the workshop is that no succinct and, more importantly, singular  
764 benchmark scenario could be envisioned that would likely elicit all or, at least, a majority of  
765 identified failures. A likely reason behind this is the lack of a comprehensive categorization of  
766 conversational failures and their triggers in mixed human-machine interactions. Having such a  
767 taxonomy would allow us to embed such triggers systematically in benchmark scenarios.

### 768 4.1 Wanted: A Taxonomy of Conversational Failures in HRI

769 Honig and Oron-Gilad (2018) recently proposed a taxonomy for failures in HRI based on a  
770 literature review of prior failure-related HRI studies. Their survey indicated a great asymmetry in  
771 these investigations, in that the majority of previous work focused on technical failures of the robot.  
772 In contrast, Honig & Oron-Gilad noticed that no strategies had been proposed to deal with “human  
773 errors”. From a conversation analytic viewpoint, the dichotomy of technical vs. human error may not



774 always be as absolute when applied to conversational failures, especially since, despite sharing some  
775 terminology, CA conceptualizes conversational success and failure quite differently. Conversation  
776 analysts conceive of successful conversation as the achievement of joint action by any party (robot or  
777 human). In this sense, when a failure occurs, the ‘blame’ lies with all participants. Similarly, success  
778 in CA terms might mean that a joint action is ‘successfully’ achieved interactionally, even if there  
779 are informational errors. For example, an invitation to meet under the clock at Grand Central station,  
780 where the recipient misunderstands the time/place might be ‘successfully’ achieved as an orderly  
781 interaction, the error being marked. In HRI, however, this failure of the ‘Schelling game’ would  
782 be considered a classic ‘grounding error’ Clark (1996), and it would certainly matter who made  
783 the error: the human or robot. While not assigning blame for some singular failure simultaneously  
784 to both participants, Uchida et al. (2019a) recently used a blame assignment strategy where the  
785 responsibility for a sequence of failures was attributed in an alternating fashion to the robot and  
786 the human. As indicated by our struggle to find a good general characterisation of conversational  
787 failures during the workshop, we advocate the construction of a taxonomy of conversational failures  
788 for mixed, that is human-machine dyads and groups. To build such a taxonomy, an interdisciplinary  
789 effort is needed, given that the types of relevant failures span the entire spectrum from the very  
790 technical (e.g. ASR errors) to the very “relational” (e.g. misunderstanding based on lack of common  
791 ground). The relevant disciplines would include linguistics, conversation analysis, robotics, NLP,  
792 HRI, and HCI. This workshop represented the first stepping stone towards this interdisciplinary  
793 effort. One theory-related advantage of taxonomy building is that it forces us to reconsider theoretical  
794 constructs from different disciplines, thereby potentially exposing gaps in the respective theories -  
795 similarly to how conversation analysis has exposed shortcomings of speech act theory (cf. Levinson,  
796 1983).

797 The process of defining the types of errors could also help us to understand why they arise, measure  
798 their impact and explore possibilities and appropriate ways to detect, mitigate and recover from  
799 them. If, for example, artificial agents and human users are mismatched conversational partners as  
800 suggested by Moore (2007) and Förster et al. (2019), and if this mismatch creates constraints and a  
801 “habitability gap” in HRI (Moore, 2017), are their specific types of failures that only occur due to  
802 such asymmetric setups? And, if yes, what does that mean for potential error management in HRI?  
803 If priors shared between interlocutors matter (Moore, 2022; Huang and Moore, 2022), how does  
804 the aligning of interactive affordances help to increase the system’s capacity to deal with errors?  
805 Moreover, errors can affect people’s perception of a robot’s trustworthiness and reliability (e.g.,  
806 Washburn et al., 2020b), as well as their acceptance and willingness to cooperate in HRI (e.g., Salem  
807 et al., 2015). What type of errors matters more? In terms of error recovery, it has been shown that  
808 social signals, such as facial action unit (AU), can enhance error detection (Stiber et al., 2023);  
809 Users’ cooperative intention can be elicited to avoid or repair from dialogue breakdowns (Uchida  
810 et al., 2019b). The question is, when facing different errors, do these strategies need to be adaptable

811 to tasks/scenarios, and if so, to what degree? Answering the above questions requires a deeper  
812 understanding of conversational failures, and taxonomy building is one possible way to increase our  
813 understanding.

814 A more practical advantage of having such a taxonomy is discussed in the next section.

## 815 4.2 Benchmarking Multimodal Speech Interfaces

816 One of the intended aims of the workshop was to define, or at least outline, some benchmark  
817 scenario that would have the “built-in” capacity to expose, if not all, at least a good number of  
818 potential communicative failures of some given speech interface. During the workshop, it became  
819 apparent that we would fail to come up with such a single scenario. It questionable whether such a  
820 scenario could exist or whether a number of scenarios would be needed to target different settings in  
821 which the speech interface is to be deployed. One main reason for our struggle that emerged during  
822 the World Café session was the lack of a taxonomy of communicative failures in HRI. Having such  
823 a taxonomy would allow the designer, or user, of a speech interface to systematically check whether  
824 it could handle the type of situation in which the identified failures are likely to occur prior to testing  
825 it “in the wild”.

826 Related to the construction of a potential (set of) benchmarks is the question of how to evaluate  
827 multimodal speech interfaces. The popular evaluation framework PARADISE Walker et al. (1997),  
828 originally designed for the assessment of unimodal dialogue systems, has already been used in  
829 multimodal HRI studies (e.g. Giuliani et al., 2013; Hwang et al., 2020; Peltason et al., 2012). Also  
830 within the HCI community multimodal alternatives to PARADISE have been proposed (e.g. Kühnel,  
831 2012). Given these existing evaluation frameworks for multimodal dialogue systems, what would a  
832 failure-based method bring to the table?

833 A characteristic of PARADISE and related frameworks is that they tend to evaluate a past dialogue  
834 according to a set of positive performance criteria. PARADISE, for example, uses measurements of  
835 *task success*, *dialogue efficiency*, and *dialogue quality* to score a given dialogue. There is likely an  
836 inverse relationship between a failure-based evaluation and, for example, *dialogue efficiency* as a  
837 dialogue containing more failures, will likely require more turns to accomplish the same task due  
838 to repair-related turns. This would mean that the efficiency of this failure-laden dialogue would be  
839 reduced. However, despite this relationship, the two methods are not commensurate. A failure-based  
840 scoring method could, for example, put positive value on the resilience of some speech interface,  
841 by assigning positive values to the number of successful repairs. This would, in some sense, be  
842 diametrically juxtaposed to efficiency measures. On the other hand, these two ways of assessing a  
843 speech interface are not mutually exclusive and could be applied simultaneously.

844 One interesting observation with respect to the surveyed studies points to a potential limitation  
845 of existing evaluation frameworks such as PARADISE. All of the referenced studies are based  
846 on turn-based interaction formats. While turn-based interaction is certainly a common format in

847 many forms of human-human and human-robot interaction, it is likely not the only one. Physical  
848 human-robot collaboration tasks which require participants to coordinate their actions in a near-  
849 simultaneous manner, for example when carrying some heavy object together, do not necessarily  
850 follow a turn-based format. While some of the involved communication channels such as speech  
851 will likely be turn-based, other channels such as sensorimotor communication (SMC, cf. Pezzulo  
852 et al., 2019) may or may not follow this format.

## 5 CONCLUSION

853 The first workshop on “Working with Troubles and Failures in Conversation between Humans and  
854 Robots” was the first effort to gather an interdisciplinary team of researchers interested in openly  
855 discuss the challenges and opportunities in designing and deploying speech interfaces for robots.  
856 Thanks to insights from conversation analysis, cognitive science, linguistics, robotics, human-robot  
857 interaction, and dialogue systems, we initiated a discussion that does not simply dismiss failures in  
858 conversational interaction as a negative outcome of the robotic system, but engages with the nature of  
859 such failures and the opportunities that arise from using them to improve the interactions. We believe  
860 this initial push will spawn a deeper research effort towards the identification of a benchmark for  
861 multimodal speech interfaces and the creation of a systematic taxonomy of failures in conversation  
862 between humans and robots which could be useful to interaction designers, both in robotics and  
863 non-robotics fields.

## 6 NOMENCLATURE

864 **Voice interfaces:** User interfaces that allow interaction with technology through spoken commands  
865 or queries.

866 **Robotic speech interfaces:** Voice interfaces applied on robots that use both speech recognition as  
867 well as synthesised or artificial voices to communicate and interact with users.

868 **Chatbots:** Text-based interfaces able to provide information, answer questions, or assist with various  
869 tasks.

870 **Agents, artificial agents, conversational agents:** Terms used interchangeably for systems designed  
871 to engage in natural language conversations with humans, by employing natural language processing  
872 and machine learning to understand and respond to user queries, provide information or assistance.

## CONFLICT OF INTEREST STATEMENT

873 Author Ioannis Papaioannou is employed by Alana AI. The remaining authors declare that the  
874 research was conducted in the absence of any commercial or financial relationships that could be  
875 construed as a potential conflict of interest.

---

## AUTHOR CONTRIBUTIONS

876 FF, MR, PH, LW, CD, JEF have organised the workshop, the contributions and notes of which form  
877 the basis of this article. FF is the lead author and has provided the main structure of the article as  
878 well as large parts of the discussion section, parts of the methods section, and overall proof-reading.  
879 MR has contributed substantial parts of the methods section, the conclusion, as well as overall  
880 proof-reading and improvements. PH, and JEF have contributed to parts of the methods section as  
881 well as overall proof-reading and improvements. FFL, SK, JH, BN, DHG, DK, JW, EEÖ, PB, GB,  
882 DP, SC, MW, LT, MP, MG, GS, PGTH, IP, DG, SA, GH, VM, EK have contributed subsections in  
883 the results section and have contributed to overall proof-reading.

## FUNDING

884 The workshop, the outcomes of which are described in this paper, was funded by the UK Engineering  
885 and Physical Sciences Research Council (EPSRC) Robotics & Autonomous Systems Network (UK-  
886 RAS) Pump Priming programme under the project title ‘Charting Current Limits and Developing  
887 Future Directions of Speech Interfaces for Robotics’.  
888 DG is supported under the EPSRC projects NLG for low-resource domains (EP/T024917/1) and  
889 CiViL (EP/T014598/1). Some of the authors are supported by the Engineering and Physical Sciences  
890 Research Council [grant number EP/V00784X/1, EP/X009343/1, EP/T014598/1] including through  
891 the Trustworthy Autonomous Systems (TAS) Hub.  
892 One of the authors has been supported by the H2020 EU project CANOPIES - A Collaborative  
893 Paradigm for Human Workers and Multi-Robot Teams in Precision Agriculture Systems, Grant  
894 Agreement 101016906.  
895 DK is funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation)  
896 under Germany’s Excellence Strategy – EXC 2002/1 “Science of Intelligence” – project number  
897 390523135.

## DATA AVAILABILITY STATEMENT

898 The original contributions presented in the study are included in the article/supplementary material,  
899 further inquiries can be directed to the corresponding author.

## REFERENCES

900 Albert, S. and Hamann, M. (2021). Putting wake words to bed: We speak wake words with  
901 systematically varied prosody, but CUIs don’t listen. In *CUI 2021 - 3rd Conference on*

- 902 *Conversational User Interfaces* (New York, NY, USA: Association for Computing Machinery),  
903 CUI '21, 1–5. doi:10.1145/3469595.3469608
- 904 Albert, S. and Ruiter, J. P. d. (2018). Repair: The Interface Between Interaction and Cognition.  
905 *Topics in Cognitive Science* 10, 279–313. doi:10.1111/tops.12339
- 906 Alharbi, S., Alrazgan, M., Alrashed, A., Alnomasi, T., Almojel, R., Alharbi, R., et al. (2021).  
907 Automatic speech recognition: Systematic literature review. *IEEE Access* 9, 131858–131876.  
908 doi:10.1109/ACCESS.2021.3112535
- 909 Axelsson, A. and Skantze, G. (2022). Multimodal user feedback during adaptive robot-human  
910 presentations. *Frontiers in Computer Science* , 135
- 911 Axelsson, A. and Skantze, G. (2023). Do you follow? a fully automated system for adaptive robot  
912 presenters. In *Proceedings of the 2023 ACM/IEEE International Conference on Human-Robot*  
913 *Interaction*. 102–111
- 914 Axelsson, N. and Skantze, G. (2020). Using knowledge graphs and behaviour trees for feedback-  
915 aware presentation agents. In *Proceedings of the 20th ACM International Conference on Intelligent*  
916 *Virtual Agents*. 1–8
- 917 Bennett, C. L., Brady, E., and Branham, S. M. (2018). Interdependence as a Frame for Assistive  
918 Technology Research and Design. In *Proceedings of the 20th International ACM SIGACCESS*  
919 *Conference on Computers and Accessibility* (New York, NY, USA: Association for Computing  
920 Machinery), ASSETS '18, 161–173. doi:10.1145/3234695.3236348
- 921 Cahya, D. E., Ramakrishnan, R., and Giuliani, M. (2019). Static and temporal differences in social  
922 signals between error-free and erroneous situations in human-robot collaboration. In *Social*  
923 *Robotics: 11th International Conference, ICSR 2019, Madrid, Spain, November 26–29, 2019,*  
924 *Proceedings 11* (Springer), 189–199
- 925 Clark, H. (1996). *Using language* (Cambridge, UK: Cambridge University Press)
- 926 Colman, M. and Healey, P. (2011). The distribution of repair in dialogue. In *Proceedings of the*  
927 *Annual Meeting of the Cognitive Science Society*. vol. 33, 1563–1568
- 928 Cuadra, A., Li, S., Lee, H., Cho, J., and Ju, W. (2021). My bad! repairing intelligent voice assistant  
929 errors improves interaction. *Proc. ACM Hum.-Comput. Interact.* 5. doi:10.1145/3449101
- 930 Curry, A. C., Papaioannou, I., Suglia, A., Agarwal, S., Shalymov, I., Xu, X., et al. (2018). Alana  
931 v2: Entertaining and informative open-domain social dialogue using ontologies and entity linking.  
932 In *1st Proceedings of Alexa Prize (Alexa Prize 2018)*
- 933 Dingemanse, M., Roberts, S. G., Baranova, J., Blythe, J., Drew, P., Floyd, S., et al. (2015). Universal  
934 principles in the repair of communication problems. *PLoS one* 10, e0136100
- 935 Dingemanse, M., Torreira, F., and Enfield, N. J. (2013). Is “Huh?” a Universal Word? Conversational  
936 Infrastructure and the Convergent Evolution of Linguistic Items. *PLoS ONE* 8, e78273. doi:10.  
937 1371/journal.pone.0078273
- 938 Dowden, B. H. (2019). *Logical Reasoning* (LibreTexts)

- 939 Enfield, N. (2017). *How We Talk: The Inner Workings of Conversation* (Hachette UK)
- 940 Enge, E. (2019). Rating the smarts of the digital personal assistants  
941 in 2018. [https://blogs.perficient.com/2018/05/01/  
942 2018-digital-personal-assistants-study/](https://blogs.perficient.com/2018/05/01/2018-digital-personal-assistants-study/), last accessed 14 June 2023
- 943 Fischer, J. E., Reeves, S., Porcheron, M., and Sikveland, R. O. (2019). Progressivity for voice  
944 interface design. In *Proceedings of the 1st International Conference on Conversational User  
945 Interfaces* (New York, NY, USA: Association for Computing Machinery), CUI '19. doi:10.1145/  
946 3342775.3342788
- 947 Fomunyam, K. G. (2020). Chaos engineering (principles of chaos engineering) as the pathway to  
948 excellence and relevance in engineering education in africa. *International Journal of Engineering  
949 and Advanced Technology (IJEAT)* 10, 146–151. doi:10.35940/ijeat.B3266.1010120
- 950 Förster, F., Saunders, J., Lehmann, H., and Nehaniv, C. L. (2019). Robots learning to say “no”:  
951 Prohibition and rejective mechanisms in acquisition of linguistic negation. *ACM Transactions on  
952 Human-Robot Interaction* 8. doi:10.1145/3359618
- 953 Foster, M. E., Craenen, B., Deshmukh, A. A., Lemon, O., Bastianelli, E., Dondrup, C., et al. (2019).  
954 Mummer: Socially intelligent human-robot interaction in public spaces. *ArXiv* abs/1909.06749
- 955 Giuliani, M., Mirnig, N., Stollnberger, G., Stadler, S., Buchner, R., and Tscheligi, M. (2015).  
956 Systematic analysis of video data from different human–robot interaction studies: a categorization  
957 of social signals during error situations. *Frontiers in Psychology* 6. doi:10.3389/fpsyg.2015.00931
- 958 Giuliani, M., Petrick, R. P., Foster, M. E., Gaschler, A., Isard, A., Pateraki, M., et al. (2013).  
959 Comparing task-based and socially intelligent behaviour in a robot bartender. In *Proceedings  
960 of the 15th ACM on International Conference on Multimodal Interaction* (New York, NY, USA:  
961 Association for Computing Machinery), ICMI '13, 263–270. doi:10.1145/2522848.2522869
- 962 Goodwin, C. (2004). A Competent Speaker Who Can't Speak: The Social Life of Aphasia. *Journal  
963 of Linguistic Anthropology* 14, 151–170. Publisher: [American Anthropological Association,  
964 Wiley]
- 965 Goodwin, C. (2007). Interactive footing. In *Reporting Talk*, eds. E. Holt and R. Clift (Cambridge:  
966 Cambridge University Press), Studies in Interactional Sociolinguistics. 16–46. doi:10.1017/  
967 CBO9780511486654.003
- 968 Green, H. N., Islam, M. M., Ali, S., and Iqbal, T. (2022). Who's laughing nao? examining perceptions  
969 of failure in a humorous robot partner. In *2022 17th ACM/IEEE International Conference on  
970 Human-Robot Interaction (HRI)*. 313–322. doi:10.1109/HRI53351.2022.9889353
- 971 Gunson, N., Hernández García, D., Sieińska, W., Dondrup, C., and Lemon, O. (2022). Developing  
972 a social conversational robot for the hospital waiting room. In *2022 31st IEEE International  
973 Conference on Robot and Human Interactive Communication (RO-MAN)*. 1352–1357. doi:10.  
974 1109/RO-MAN53752.2022.9900827

- 975 Hancock, P. A., Billings, D. R., Schaefer, K. E., Chen, J. Y. C., de Visser, E. J., and Parasuraman, R.  
976 (2011). A meta-analysis of factors affecting trust in human-robot interaction. *Human Factors* 53,  
977 517–527. doi:10.1177/0018720811417254. PMID: 22046724
- 978 Healey, P. (2008). Interactive misalignment: The role of repair in the development of group  
979 sub-languages. *Language in Flux. College Publications* 212
- 980 Healey, P., Plant, N., Howes, C., and Lavelle, M. (2015). When words fail: Collaborative gestures  
981 during clarification dialogues. In *2015 AAAI Spring Symposium Series*
- 982 Healey, P. G. (1997). Expertise or expertese?: The emergence of task-oriented sub-languages. In  
983 *Proceedings of the 19th annual conference of the cognitive science society* (Stanford University  
984 Stanford, CA), 301–306
- 985 Healey, P. G., De Ruiter, J. P., and Mills, G. J. (2018a). Editors' introduction: miscommunication.  
986 *Topics in Cognitive Science* 10, 264–278
- 987 Healey, P. G., Mills, G. J., Eshghi, A., and Howes, C. (2018b). Running repairs: Coordinating  
988 meaning in dialogue. *Topics in cognitive science* 10, 367–388
- 989 Heldner, M. and Edlund, J. (2010). Pauses, gaps and overlaps in conversations. *Journal of Phonetics*  
990 38, 555–568
- 991 Hirayama, N., Yoshino, K., Itoyama, K., Mori, S., and Okuno, H. G. (2015). Automatic  
992 speech recognition for mixed dialect utterances by mixing dialect language models. *IEEE/ACM*  
993 *Transactions on Audio, Speech, and Language Processing* 23, 373–382
- 994 Holtgraves, T. M. (2013). *Language as social action: Social psychology and language use*  
995 (Psychology Press)
- 996 Holthaus, P., Schulz, T., Lakatos, G., and Soma, R. (2023). Communicative Robot Signals:  
997 Presenting a New Typology for Human-Robot Interaction. In *International Conference on*  
998 *Human-Robot Interaction (HRI 2023)* (Stockholm, Sweden: ACM/IEEE), 132–141. doi:10.1145/  
999 3568162.3578631
- 1000 Honig, S. and Oron-Gilad, T. (2018). Understanding and resolving failures in human-robot  
1001 interaction: Literature review and model development. *Frontiers in Psychology* 9. doi:10.  
1002 3389/fpsyg.2018.00861
- 1003 Hough, J., de Kok, I., Schlangen, D., and Kopp, S. (2015). Timing and grounding in motor skill  
1004 coaching interaction: Consequences for the information state. In *Proceedings of the 19th SemDial*  
1005 *Workshop on the Semantics and Pragmatics of Dialogue (goDIAL)*. 86–94
- 1006 Hough, J. and Schlangen, D. (2016). Investigating fluidity for human-robot interaction with real-  
1007 time, real-world grounding strategies. In *Proceedings of the 17th Annual Meeting of the Special*  
1008 *Interest Group on Discourse and Dialogue* (Los Angeles: ACL), 288–298
- 1009 Howes, C. and Eshghi, A. (2021). Feedback relevance spaces: Interactional constraints on processing  
1010 contexts in dynamic syntax. *Journal of Logic, Language and Information* 30, 331–362

- 1011 Howes, C., Lavelle, M., Healey, P., Hough, J., and McCabe, R. (2016). Helping hands? gesture  
1012 and self-repair in schizophrenia. In *Proceedings of the Resources and Processing of Linguistic  
1013 and Extra-Linguistic Data from People with Various Forms of Cognitive/Psychiatric Impairments  
1014 (RaPID-2016)*. 9–13
- 1015 Huang, G. and Moore, R. K. (2022). Is honesty the best policy for mismatched partners? aligning  
1016 multi-modal affordances of a social robot: an opinion paper. *Frontiers in Virtual Reality*
- 1017 Hwang, E. J., Kyu Ahn, B., Macdonald, B. A., and Seok Ahn, H. (2020). Demonstration of hospital  
1018 receptionist robot with extended hybrid code network to select responses and gestures. In *2020  
1019 IEEE International Conference on Robotics and Automation (ICRA)*. 8013–8018. doi:10.1109/  
1020 ICRA40945.2020.9197160
- 1021 Kendrick, K. H. and Drew, P. (2016). Recruitment: Offers, Requests, and the  
1022 Organization of Assistance in Interaction. *Research on Language and Social Interaction*  
1023 49, 1–19. doi:10.1080/08351813.2016.1126436. Publisher: Routledge \_eprint:  
1024 <https://doi.org/10.1080/08351813.2016.1126436>
- 1025 Kontogiorgos, D., Pereira, A., Sahindal, B., van Waveren, S., and Gustafson, J. (2020a). Behavioural  
1026 responses to robot conversational failures. In *Proceedings of the 2020 ACM/IEEE International  
1027 Conference on Human-Robot Interaction*. 53–62
- 1028 Kontogiorgos, D., Tran, M., Gustafson, J., and Soleymani, M. (2021). A systematic cross-  
1029 corpus analysis of human reactions to robot conversational failures. In *Proceedings of the  
1030 2021 International Conference on Multimodal Interaction*. 112–120
- 1031 Kontogiorgos, D., Van Waveren, S., Wallberg, O., Pereira, A., Leite, I., and Gustafson, J. (2020b).  
1032 Embodiment effects in interactions with failing robots. In *Proceedings of the 2020 CHI conference  
1033 on human factors in computing systems*. 1–14
- 1034 Krauss, R. M. and Fussell, S. R. (1990). Mutual knowledge and communicative effectiveness.  
1035 *Intellectual teamwork: Social and technological foundations of cooperative work*, 111–146
- 1036 Kühnel, C. (2012). *Quantifying Quality Aspects of Multimodal Interactive Systems* (Springer Science  
1037 & Business Media)
- 1038 Kwon, M., Huang, S. H., and Dragan, A. D. (2018). Expressing robot incapability. In *Proceedings  
1039 of the 2018 ACM/IEEE International Conference on Human-Robot Interaction*. 87–95
- 1040 Lau, J., Zimmerman, B., and Schaub, F. (2018). Alexa, are you listening? privacy perceptions,  
1041 concerns and privacy-seeking behaviors with smart speakers. *Proceedings of the ACM on  
1042 Human-Computer Interaction* 2, 1–31
- 1043 Levinson, S. C. (1983). *Pragmatics* (Cambridge, UK: Cambridge University Press)
- 1044 Liza, F. F. (2019). *Improving Training of Deep Neural Network Sequence Models* (University of  
1045 Kent (United Kingdom))
- 1046 Luong, H.-T. and Yamagishi, J. (2020). Nautilus: a versatile voice cloning system. *IEEE/ACM  
1047 Transactions on Audio, Speech, and Language Processing* 28, 2967–2981



- 1048 Maraev, V., Bernardy, J.-P., and Howes, C. (2021). Non-humorous use of laughter in spoken dialogue  
1049 systems. In *Linguistic and Cognitive Approaches to Dialog Agents (LaCATODA 2021)*. 33–44
- 1050 Marge, M., Espy-Wilson, C., Ward, N. G., Alwan, A., Artzi, Y., Bansal, M., et al. (2022). Spoken  
1051 language interaction with robots: Recommendations for future research. *Computer Speech &*  
1052 *Language* 71, 101255. doi:<https://doi.org/10.1016/j.csl.2021.101255>
- 1053 Mirnig, N., Giuliani, M., Stollnberger, G., Stadler, S., Buchner, R., and Tscheligi, M. (2015). Impact  
1054 of robot actions on social signals and reaction times in hri error situations. In *Social Robotics:*  
1055 *7th International Conference, ICSR 2015, Paris, France, October 26-30, 2015, Proceedings 7*  
1056 (Springer), 461–471
- 1057 Mirnig, N., Stollnberger, G., Miksch, M., Stadler, S., Giuliani, M., and Tscheligi, M. (2017). To err  
1058 is robot: How humans assess and act toward an erroneous social robot. *Frontiers in Robotics and*  
1059 *AI* , 21
- 1060 Moore, R. K. (2007). Spoken language processing: Piecing together the puzzle. *Speech*  
1061 *communication* 49, 418–435
- 1062 Moore, R. K. (2017). Is spoken language all-or-nothing? implications for future speech-based human-  
1063 machine interaction. *Dialogues with Social Robots: Enablements, Analyses, and Evaluation* ,  
1064 281–291
- 1065 Moore, R. K. (2022). Whither the priors for (vocal) interactivity? *arXiv preprint arXiv:2203.08578*
- 1066 Nielsen, J. (1995). 10 usability heuristics for user interface design. <https://www.nngroup.com/articles/ten-usability-heuristics/>, last accessed 14 June 2023
- 1068 Ozkan, E. E., Gurion, T., Hough, J., Healey, P. G., and Jamone, L. (2022). Speaker motion patterns  
1069 during self-repairs in natural dialogue. In *Companion Publication of the 2022 International*  
1070 *Conference on Multimodal Interaction* (New York, NY, USA: Association for Computing  
1071 Machinery), ICMI '22 Companion, 24–29. doi:10.1145/3536220.3563684
- 1072 Papaioannou, I., Cercas Curry, A., Part, J. L., Shalymov, I., Xu, X., Yu, Y., et al. (2017). Alana:  
1073 Social dialogue using an ensemble model and a ranker trained on user feedback. *Proc. AWS re:*  
1074 *INVENT*
- 1075 Park, S., Healey, P. G. T., and Kaniadakis, A. (2021). Should robots blush? In *Proceedings*  
1076 *of the 2021 CHI Conference on Human Factors in Computing Systems* (New York, NY, USA:  
1077 Association for Computing Machinery), CHI '21. doi:10.1145/3411764.3445561
- 1078 Peltason, J., Riether, N., Wrede, B., and Lütkebohle, I. (2012). Talking with robots about objects: A  
1079 system-level evaluation in hri. In *Proceedings of the Seventh Annual ACM/IEEE International*  
1080 *Conference on Human-Robot Interaction* (New York, NY, USA: Association for Computing  
1081 Machinery), HRI '12, 479–486. doi:10.1145/2157689.2157841
- 1082 Pezzulo, G., Donnarumma, F., Dindo, H., D'Ausilio, A., Konvalinka, I., and Castelfranchi, C. (2019).  
1083 The body talks: Sensorimotor communication and its brain and kinematic signatures. *Physics of*  
1084 *Life Reviews* 28, 1–21. doi:<https://doi.org/10.1016/j.plrev.2018.06.014>

- 1085 Porcheron, M., Fischer, J. E., Reeves, S., and Sharples, S. (2018). Voice interfaces in everyday life.  
1086 In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems* (New York,  
1087 NY, USA: Association for Computing Machinery), CHI '18, 1–12. doi:10.1145/3173574.3174214
- 1088 Purver, M. (2004). *The Theory and Use of Clarification Requests in Dialogue*. Ph.D. thesis, King's  
1089 College, University of London
- 1090 Purver, M., Eshghi, A., and Hough, J. (2011). Incremental semantic construction in a dialogue  
1091 system. In *Proceedings of the ninth international conference on computational semantics (IWCS*  
1092 *2011)*. 365–369
- 1093 Ragni, M., Rudenko, A., Kuhnert, B., and Arras, K. O. (2016). Errare humanum est: Erroneous  
1094 robots in human-robot interaction. In *2016 25th IEEE International Symposium on Robot and*  
1095 *Human Interactive Communication (RO-MAN)*. 501–506. doi:10.1109/ROMAN.2016.7745164
- 1096 Ram, A., Prasad, R., Khatri, C., Venkatesh, A., Gabriel, R., Liu, Q., et al. (2018). Conversational ai:  
1097 The science behind the alexa prize. *arXiv preprint arXiv:1801.03604*
- 1098 Salem, M., Lakatos, G., Amirabdollahian, F., and Dautenhahn, K. (2015). Would you trust a  
1099 (faulty) robot? effects of error, task type and personality on human-robot cooperation and trust. In  
1100 *Proceedings of the tenth annual ACM/IEEE international conference on human-robot interaction*.  
1101 141–148
- 1102 Schegloff, E. A. (1992a). Repair after next turn: The last structurally provided defense of  
1103 intersubjectivity in conversation. *American journal of sociology* 97, 1295–1345
- 1104 Schegloff, E. A. (1992b). Repair after next turn: The last structurally provided defense of  
1105 intersubjectivity in conversation. *American Journal of Sociology* 97, 1295–1345
- 1106 Schegloff, E. A. (1997). Third turn repair. *AMSTERDAM STUDIES IN THE THEORY AND*  
1107 *HISTORY OF LINGUISTIC SCIENCE SERIES 4*, 31–40
- 1108 Schegloff, E. A. (2007). *Sequence organization in interaction: Volume 1: A primer in conversation*  
1109 *analysis* (New York: Cambridge University Press)
- 1110 Schegloff, E. A., Jefferson, G., and Sacks, H. (1977a). The preference for self-correction in the  
1111 organization of repair in conversation. *Language* 53, 361–382
- 1112 Schegloff, E. A., Jefferson, G. D., and Sacks, H. (1977b). The preference for self-correction in the  
1113 organization of repair in conversation. *Language* 53, 361 – 382
- 1114 Shriberg, E. (1994). *Preliminaries to a Theory of Speech Disfluencies*. Ph.D. thesis, University of  
1115 California, Berkeley
- 1116 Siwach, G., Haridas, A., and Chinni, N. (2022). Evaluating operational readiness using chaos  
1117 engineering simulations on kubernetes architecture in big data. In *2022 International Conference*  
1118 *on Smart Applications, Communications and Networking (SmartNets)*. 1–7. doi:10.1109/  
1119 SmartNets55823.2022.9993998
- 1120 Skantze, G. (2005). Exploring human error recovery strategies: Implications for spoken dialogue  
1121 systems. *Speech Communication* 45, 325–341

- 1122 Skantze, G. and Doğruöz, A. S. (2023). The open-domain paradox for chatbots: Common ground as  
1123 the basis for human-like dialogue. *arXiv preprint arXiv:2303.11708*
- 1124 Stiber, M., Taylor, R. H., and Huang, C.-M. (2023). On using social signals to enable flexible error-  
1125 aware hri. In *Proceedings of the 2023 ACM/IEEE International Conference on Human-Robot*  
1126 *Interaction* (New York, NY, USA: Association for Computing Machinery), HRI '23, 222–230.  
1127 doi:10.1145/3568162.3576990
- 1128 Strathearn, C. and Gkatzia, D. (2021a). Chefbot: A novel framework for the generation of  
1129 commonsense-enhanced responses for task-based dialogue systems. In *Proceedings of the 14th*  
1130 *International Conference on Natural Language Generation* (Aberdeen, Scotland, UK: Association  
1131 for Computational Linguistics), 46–47
- 1132 Strathearn, C. and Gkatzia, D. (2021b). Task2Dial dataset: A novel dataset for commonsense-  
1133 enhanced task-based dialogue grounded in documents. In *Proceedings of the 4th International*  
1134 *Conference on Natural Language and Speech Processing (ICNLSP 2021)* (Trento, Italy:  
1135 Association for Computational Linguistics), 242–251
- 1136 Tian, L. and Oviatt, S. (2021). A Taxonomy of Social Errors in Human-Robot Interaction. *ACM*  
1137 *Transactions on Human-Robot Interaction (THRI)* 10, 1–32
- 1138 Tomasello, M. (2009). *Why we cooperate* (MIT press)
- 1139 Trung, P., Giuliani, M., Miksch, M., Stollnberger, G., Stadler, S., Mirnig, N., et al. (2017). Head  
1140 and shoulders: Automatic error detection in human-robot interaction. In *Proceedings of the 19th*  
1141 *ACM International Conference on Multimodal Interaction* (New York, NY, USA: Association for  
1142 Computing Machinery), ICMI '17, 181–188. doi:10.1145/3136755.3136785
- 1143 Uchida, T., Minato, T., Koyama, T., and Ishiguro, H. (2019a). Who is responsible for a dialogue  
1144 breakdown? an error recovery strategy that promotes cooperative intentions from humans by  
1145 mutual attribution of responsibility in human-robot dialogues. *Frontiers in Robotics and AI* 6,  
1146 doi:10.3389/frobt.2019.00029
- 1147 Uchida, T., Minato, T., Koyama, T., and Ishiguro, H. (2019b). Who is responsible for a dialogue  
1148 breakdown? an error recovery strategy that promotes cooperative intentions from humans by  
1149 mutual attribution of responsibility in human-robot dialogues. *Frontiers in Robotics and AI* 6, 29
- 1150 Walker, M. A., Litman, D. J., Kamm, C. A., and Abella, A. (1997). Paradise: A framework for  
1151 evaluating spoken dialogue agents. In *Proceedings of the 35th Annual Meeting of the Association*  
1152 *for Computational Linguistics and Eighth Conference of the European Chapter of the Association*  
1153 *for Computational Linguistics* (USA: Association for Computational Linguistics), ACL '98/EACL  
1154 '98, 271–280. doi:10.3115/976909.979652
- 1155 Washburn, A., Adeleye, A., An, T., and Riek, L. D. (2020a). Robot errors in proximate hri:  
1156 How functionality framing affects perceived reliability and trust. *J. Hum.-Robot Interact.* 9,  
1157 doi:10.1145/3380783

- 1158 Washburn, A., Adeleye, A., An, T., and Riek, L. D. (2020b). Robot errors in proximate hri: how  
1159 functionality framing affects perceived reliability and trust. *ACM Transactions on Human-Robot*  
1160 *Interaction (THRI)* 9, 1–21
- 1161 Williams, J., Fong, J., Cooper, E., and Yamagishi, J. (2021a). Exploring Disentanglement with  
1162 Multilingual and Monolingual VQ-VAE. In *Proc. 11th ISCA Speech Synthesis Workshop (SSW*  
1163 *11)*. 124–129. doi:10.21437/SSW.2021-22
- 1164 Williams, J., Pizzi, K., Das, S., and Noé, P.-G. (2022). New challenges for content privacy in speech  
1165 and audio. In *Proc. 2nd ISCA Symposium on Security and Privacy in Speech Communication*.  
1166 1–6. doi:10.21437/SPSC.2022-1
- 1167 Williams, J., Zhao, Y., Cooper, E., and Yamagishi, J. (2021b). Learning disentangled phone  
1168 and speaker representations in a semi-supervised vq-vae paradigm. In *ICASSP 2021-2021*  
1169 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (IEEE),  
1170 7053–7057
- 1171 Xu, H., Khassanov, Y., Zeng, Z., Chng, E. S., Ni, C., Ma, B., et al. (2020). Independent language  
1172 modeling architecture for end-to-end asr. In *ICASSP 2020-2020 IEEE International Conference*  
1173 *on Acoustics, Speech and Signal Processing (ICASSP)* (IEEE), 7059–7063
- 1174 Özkan, E. E., Gurion, T., Hough, J., Healey, P. G., and Jamone, L. (2021). Specific hand motion  
1175 patterns correlate to miscommunications during dyadic conversations. In *2021 IEEE International*  
1176 *Conference on Development and Learning (ICDL)*. 1–6. doi:10.1109/ICDL49984.2021.9515613
- 1177 Özkan, E. E., Healey, P. G., Gurion, T., Hough, J., and Jamone, L. (2023). Speakers raise their  
1178 hands and head during self-repairs in dyadic conversations. *IEEE Transactions on Cognitive and*  
1179 *Developmental Systems* , 1–1doi:10.1109/TCDS.2023.3254808