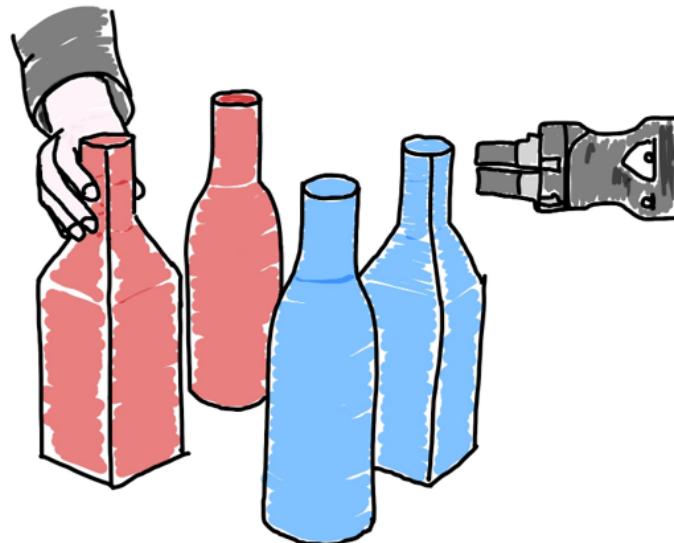


1     **Are Explanations All We Need? Investigating the Impact of XAI on Robot**  
2     **Failures and Trust Recovery in Human-Robot Cooperation Tasks**  
3

4     ANONYMOUS AUTHOR(S)\*  
5  
6



27     Fig. 1. Cooperative human-robot interaction task: Both, robot and human, sorting bottles regarding their shapes.  
28

29     Failures are part of our lives. We have learned to deal with failures and fix them in private and work-life. However, with the emerging  
30     use of robots in an industrial context, we have to learn to deal with new error originators: robots. Studies show that robot errors  
31     decrease user trust after a robot failure, negatively impacting performance. The question is whether it is possible to restore the lost  
32     trust. Based on the literature, we identify two steps that can support trust recovery: (1) communicating failures with the help of XAI  
33     and (2) fixing robot failures. Further, we illustrate a setup for a planned study that investigates both aspects and end the paper with a  
34     discussion of open questions regarding human-centered and robot-centered issues that should be considered in failure mitigation.  
35

36     CCS Concepts: • **Human-centered computing** → **Empirical studies in HCI; User centered design; User studies.**  
37

38     Additional Key Words and Phrases: industrial robots, failure mitigation, trust recovery, user-study, human-centered AI  
39

40     **ACM Reference Format:**

41     Anonymous Author(s). 2023. Are Explanations All We Need? Investigating the Impact of XAI on Robot Failures and Trust Recovery in  
42     Human-Robot Cooperation Tasks. In *ACM Conference on Human Factors in Computing Systems, April 23–28, 2023, Hamburg, Germany*.  
43     ACM, New York, NY, USA, 10 pages. <https://doi.org/XXXXXX.XXXXXXX>

44  
45     Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not  
46     made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components  
47     of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to  
48     redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

49     © 2023 Association for Computing Machinery.  
50     Manuscript submitted to ACM

## 53    1 INTRODUCTION

54  
55 In the course of our lives, we make failures. As children, we learn how to deal with them by understanding failure and  
56 its effects, looking for solutions and trying to avoid them in the future. This task is often challenging and accompanies  
57 us throughout our lives. In addition to our failures, there are also the ones made by fellow human beings that we must  
58 learn to deal with. Failures occur when people interact and work together on a task. Communication is one crucial  
59 point that helps us correct and avoids these failures in the future. We get in touch with others, report to them when an  
60 error has occurred, and our counterpart can explain how it happened or how the problem could be fixed or avoided.  
61 Especially in the work context, investigating failures is vital to ensure safety [17]. Nowadays, more and more companies  
62 use industrial robots in addition to human workers. These can also cause failures that humans have to deal with. The  
63 investigation of robot failures in human perception and behaviour is the subject of numerous studies. In particular, the  
64 subject of research is the question of how a robot failure, which is often accompanied by a loss of trust on the part of  
65 the human, can be returned to trusting cooperation [24]. One possibility is the use of explanations. The research area  
66 of Explainable AI (XAI) deals with explaining decisions or internal states of AI systems. Therefore, XAI supports a  
67 human-centered AI design meaning that AI systems take human needs into account [18]. Whether explanations help  
68 restore trust in human-robot cooperation was the subject of investigation in the work of Hald et al. [10]. The authors  
69 investigated the influence of explanations on trust recovery after a robot failure in a VR setting. The authors concluded  
70 that explanations alone could not restore trust in an industrial robot after a robot failure.  
71

72 If an explanation is insufficient to restore trust, the logical next step is realising the solution given in the explanation  
73 and, in the process, fixing the failure. This workshop paper provides an overview of research related to the influence of  
74 fixing the robot's failure on the recovery of trust. Based on existing literature, we present explanations in the context of  
75 robot failures and introduce the combination of explanations with user error correction as a possible helpful form of  
76 interaction after a robot failure. Furthermore, we reason about the implications for a planned follow-up study extending  
77 the results of Hald et al. [10]. Finally, our paper concludes with a discussion of open questions regarding human-centered  
78 and robot-centered issues in this context. Our work contributes towards investigating XAI in real-life use cases where  
79 humans and robots must successfully interact to solve tasks.  
80

## 81    2 BACKGROUND

### 82    2.1 Robot Failures

83 When robots show not-intended behaviour, terms like failure, fault, or error describe this situation [12]. A failure  
84 describes a “behavior or service being performed by the system to deviate from the ideal, normal, or correct functionality”  
85 [2, p. 9]. Failures are caused by errors in the system (e.g., mechanical errors) [12]. The causes of errors are faults (e.g., a  
86 loose screw) [12]. Honig and Oron-Gilad [12] underline that in human-robot interaction, robot failures are common  
87 due to complex and often unstructured situations.  
88

### 89    2.2 Trust

90 Authors such as de Visser et al. [5], Hancock et al. [11], Lee and See [13] point out the importance of an appropriate  
91 level of trust towards a robot, as too much trust can lead to misuse of the robot and thus dangerous situations. On the  
92 other hand, too little trust can lead to the robot not being used optimally.  
93

94 Therefore, studying factors influencing trust between humans and machines has been an important research subject  
95 for decades. Different definitions of trust can be found in human-robot interaction. One of the most commonly used  
96

definitions of interaction is that of Lee and See [13]. They define trust as “the attitude that an agent will help achieve an individual’s goals in a situation characterized by uncertainty and vulnerability.”[13, p. 54].

Hancock et al. [11] systematically summarised possible factors influencing trust. They point out three areas influencing people’s trust in robots: robot-related, environment-related, and user-related factors. Robot-related factors include the robot’s personality and behaviour and the reliability and predictability of the robot’s behaviour. Environment-related factors include the nature and complexity of the task, as well as aspects such as culture and communication. Finally, user-related factors include, for example, the expertise and competence of users.

When investigating the reliability of robots, Salem et al. [20] found that a robot is perceived as less trustworthy after a failure. However, despite the decrease in trust, participants followed the robot’s instructions.

### 3 DEALING WITH ROBOT FAILURES

When an error occurs in human-robot interaction, users often cannot understand how the failure happened, how it can be fixed and how it can be avoided in the future. Lack of understanding leads to a degradation of performance and mistrust of the user. Therefore, it is vital to address and solve the lack of understanding preventing loss of performance and trust.

Honig and Oron-Gilad [12] present their Robot Failure Human Information Processing Model. Their model addresses the mitigation of robot failure from a user-centred perspective by taking into account the user’s cognitive abilities [12, p. 8]:

- **Communicating Failures:** This is about how a robot failure is communicated. Communication can be done with visual indicators (e.g. lights), screens with additional information or audio and speech (e.g. sounds, up to verbal communication).
- **Perception & Comprehension of Failures:** This aspect deals with the user’s perception of the failure. Only when the user perceives the failure can countermeasures be initiated. Here, the user must understand and react to the robot’s failure. The user’s mental model also influences how they deal with a robot’s failure. An incorrect model of the robot’s behaviour can make it difficult to solve the problem after a failure.
- **Solving Failures:** User motivation is required to resolve a robot failure. In addition, users must decide how to fix the failure and translate this into action.

On this basis, strategies for dealing with robot failures can be developed for the robot or the user. Honig and Oron-Gilad [12] describes three different strategies: (1) set expectations (i.e., the user can correctly assess the failure potential of the robot), (2) communicate correctly (i.e., how and what is communicated in a failure situation) and (3) ask for help (i.e., the robot asks the user for help).

## 4 COMMUNICATING & FIXING ROBOT FAILURES

### 4.1 Communicating Robot Failure

In addition to harm reduction strategies, as stated before, providing explanations of the robot’s behaviour and decisions could also help to increase user trust. The field of XAI has the goal to help users to “understand, appropriately trust, and effectively manage”[9, p.44] AI partners. de Visser et al. [5] provide an overview of trust in human-robot teams. They emphasise the importance of trust for successful human-robot tasks. They mention two essential directions for (re)building human-robot relationships through trust calibration: (1) predictive and preventive trust to minimise the risk of trust violation, and (2) reactive and reparative trust when trust is violated.

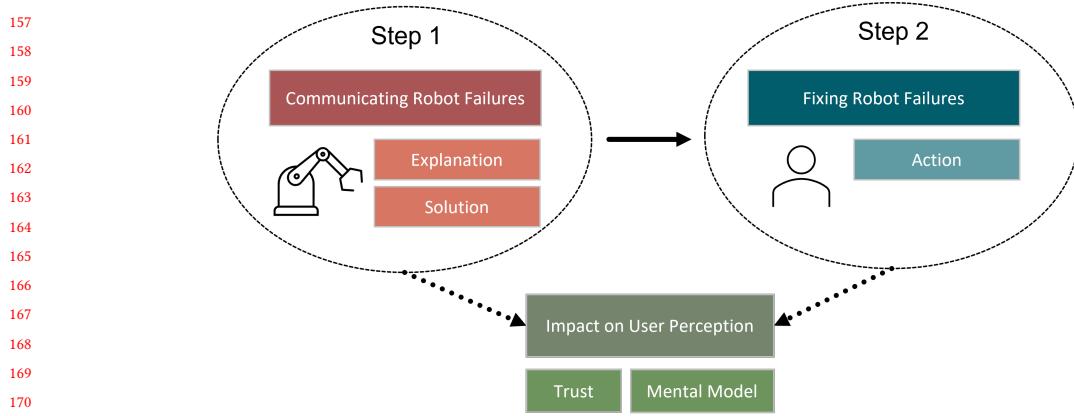


Fig. 2. Based on the literature, we conclude two steps for mitigating robot failures: first, a robot failure is communicated to the user with the help of an explanation and a possible solution. Second, the user comes into action by actually fixing the error. The steps can, but do not need to rely on each other and have an impact on users' perception of the robot.

One method for the predictive and reparative trust that de Visser et al. [5] suggests is using explanations. They suggest that explanations help calibrate trust by providing users with information about the robot. XAI approaches aim to gain insights into a robot's behaviour and goals. For example, Sheh [22] recommends using decision trees in explanatory human-robot dialogues. They argue that decision trees have two significant advantages: (1) the intrinsic explainability of if-then-else rules and (2) attribute-centric explanations that do not refer to a specific example but to attributes relevant to different actions. Therefore, they aim to use decision trees in a human-robot dialogue and enrich it with visual explanation modalities such as histograms. The work of Das et al. [4] uses handwritten explanations to study three different types of errors (i.e. navigation, arm motion planning and object recognition) of an industrial robot arm. They distinguish four different contents of explanations [4, p. 354]:

- **Action Based:** e.g., "Robot could not find the object."
- **Context Based:** e.g., "Robot could not find the object because the object is hidden from view."
- **Action Based History:** e.g., "The robot finished scanning objects at its current location, but could not find the desired object."
- **Context Based History:** e.g., "The robot finished scanning objects at its current location, but could not find the desired object because the desired object is hidden from view."

In an online user study, videos were shown where the robot performed a task. After three videos where the robot solved the task, participants saw six failure simulations without explanations. Afterwards, participants were asked to identify the cause of the failure and possible solutions. Finally, another 12 failure videos with explanations were shown. Again, participants had to identify the cause of the failure and possible solutions. The results indicate that context-based explanations are best suited for end users to identify the reason for a robot failure and find a helpful solution.

In the work of Hald et al. [10], the influence of explanations on the trust of the robot after a robot error is investigated. For this purpose, Hald et al. [10] used an interactive scenario where study participants had to sort bottles by shape with a robot. The participant had the task to sort the red bottles according to their shape, such that all red rectangular bottles will be within the white area on the left side and all round-shaped, red bottles on the right side on the white basis (see upper part in Figure 3 on page 6). The task for the cobot is identical, except its task comprises the blue bottles

209 instead of the red ones. The setting was realised in virtual reality (VR). In total three different conditions were tested  
210 for in the study:

- 211
- 212 • **No Explanation:** After the robot error occurred, participants only received information that an error had  
213 occurred and that the task had not been successfully completed.
  - 214 • **Explanation:** Participants received an explanation of why the robot error had occurred. The explanation was  
215 projected onto the virtual table in front of them, and stated: "A computer vision error occurred. The system did  
216 not successfully distinguish the shapes in the current lighting conditions." (see Figure 3 on the next page)
  - 217 • **Explanation and Solution:** In addition to the explanation, participants received information on how to avoid  
218 the error in the future.

219

220 As already found in numerous studies by other authors, Hald et al. [10] also found a loss of trust after the robot made  
221 a mistake. However, the three different conditions did not differ, i.e., no influence of explanations on a trust recovery  
222 after a robot error could be found. However, at the end of the study, the participants rated the explanations as helpful to  
223 trust or distrust the robot.

## 227 4.2 Fixing Robot Failure

228 So far, most of the research regarding trust in (faulty) robots focused on communication strategies, e.g., denying,  
229 apologizing for or explaining the failure [7], type of failure, temporal occurrence of failure [19], the severity of the  
230 failure [19] or the personal relevance of failures [8]. Also, characteristics of the robot, e.g., its apparel, have been the  
231 object of research on its influence on the trust. To the best of our knowledge, investigating trust repair strategies that  
232 involve resolving the failure after its occurrence has not been the main focus of any research activity. Morales et al.  
233 [16] investigated whether one is willing to assist a robot in case it not being able to grab an object. They combined the  
234 occurrence of the failure with other interactions that included the robot behaving in a harmful way, e.g., throwing a  
235 (foam) potato at the participant or crushing a bag of potato chips. The order of the different types of interaction had  
236 an influence on the willingness to help the robot. Specifically, the willingness decreased significantly when the risky  
237 behaviour, independent of the participant or an object potentially being harmed, occurred before the robot failed at  
238 grabbing a box [16]. Morales et al. [16] show that successfully resolving a robot's failure is affected by its behaviour.

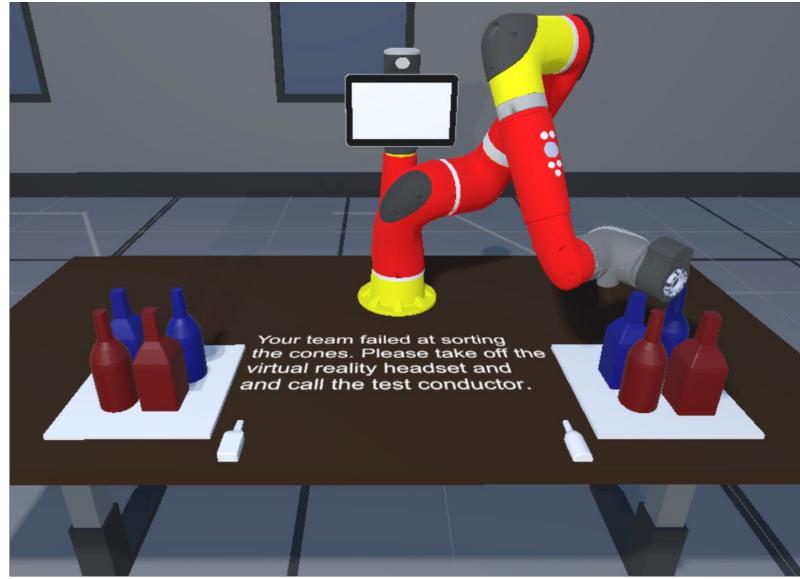
239 Explaining the reason for a failure without fixing the failure mitigates the loss of trust after the failure [10]. Adding  
240 the proposal of a possible solution, without the possibility of correcting the failure, did not gain any further benefit  
241 in mitigating the trust loss. In Hald et al. [10] resolving the failure itself was not tested for. However, the reason for  
242 the failure of the robot as well as the solution resolving the failure was given in form of a text-based explanation (see  
243 Figure Figure 3 on the following page).

244 The logical follow-up question is: What influence has fixing the failure upon the explanation and/or solution given  
245 to one? Thus, we will use the experimental setup by Hald et al. [10] as a basis and adapt it to the needs for answering  
246 our research question (see figure 4 on page 7. For this, we made three changes.

247 Firstly, we will change the VR setting to real life as Hald et al. [10] already outlined in their look in their outlook.

248 One might argue that there are significant differences between VR and reality complicating the comparison of  
249 results between those different settings. While this seems to be true for the perception of proxemics [6] and safety [15].  
250 Stressing participants with the Trier Social Stress Test elicited similar subjective ratings, EEG and EDA results when  
251 comparing the VR and real setting. Only the cortisol level differed between the conditions [23]. Also, the perception of  
252 eeriness and likeability did not differ between the presentation of a humanoid robot in real life or in VR [14]. Ünal et al.  
253

261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280



281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300



301 Fig. 3. The virtual setup (top) presented in Hald et al. [10] serves as a basis for the investigation of the helpfulness of explanations in  
302 solving robot failures in a real-life setup (bottom).

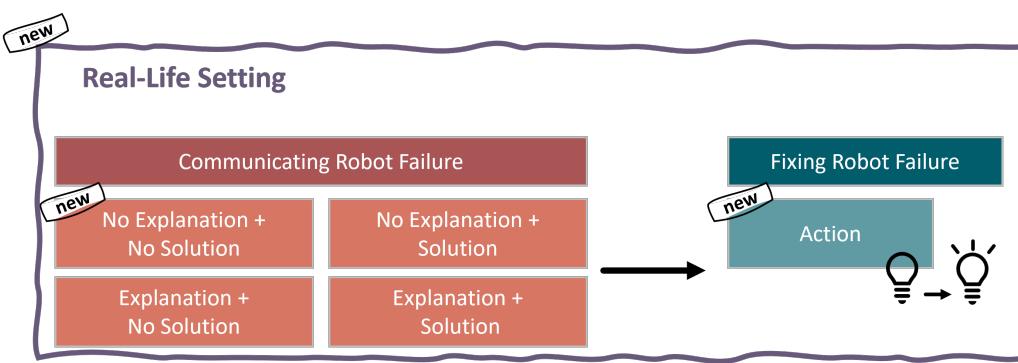
303  
304  
305  
306  
307  
308  
309  
310  
311

[26] found no difference when comparing the restorative effect of being in a real-life natural environment or a similar VR environment. We cannot rule out that our experimental setup in VR or in real life might have an influence on the results of explanations and/or solving the failure on trusting the robot or not, thus, we shall be very cautious when comparing results.

312

313 Secondly, we will extend the number of conditions. The additional fourth condition entails no explanation, but a  
 314 proposal of a solution. This allows examining both factors (no) explanation and (no) solutions separately as well as in  
 315 interaction.

316 Thirdly and most important change, the failure will be fixed in all four conditions. The **explanation and solution**  
 317 condition of the previous study entailed the information that "*better lightning conditions will help with successful sorting*"  
 318 [10, p. 223]. Consequently, a reasonable fix for the failure is turning on the light. For the sake of maintaining the  
 319 possibility of comparing all four conditions, we must have a coherent way of acting upon the proposed solution to the  
 320 failure. It would be presumptuous to expect participants in the condition without any explanation or solution to initiate  
 321 turning on the light. While, it would be of interest to leave the action to the participants, as it might give them a feeling  
 322 of greater involvement with the robot, the deviation between conditions would be too great to allow any comparison.  
 323 Thus, standardizing across conditions by leaving the action upon the study conductor to turn on the light allows a  
 324 sound comparison.



342 Fig. 4. The extended study design compared to the work of Hald et al. [10]. The real-life setting, an “explanation and no solution”  
 343 condition, and the action of the study conductor to fix the robot failure are new components of the planned study.

## 344 5 DISCUSSION

345 With the perspective of the great increase of robots advancing into several areas of our life, more and more people,  
 346 especially those without great knowledge of technical systems, will interact with robots and observe robots fail and  
 347 possibly lose trust. Thus, the communication of failures, but also the prevention of as well as the intervention after the  
 348 occurrence of failures must be further researched to allow robots to be a further supporting component in facilitating  
 349 our lives. A list of open questions related to our intended experiment, but also beyond the scope it is being touched  
 350 upon.

351 *Who is responsible for robot failure mitigation?* When dealing with failure mitigation, the question inevitably arises:  
 352 Who is responsible for fixing the failure? The person (or robot) who caused the failure? Or the person who has an idea  
 353 for a solution? Or the person (or robot) who can implement the solution? Often, humans have to correct a machine’s  
 354 failures because people are good at recognising and solving failures, or at least not making them worse, based on the  
 355 situation’s circumstances and the available options. In communication situations, it becomes apparent that people are  
 356 primarily responsible for resolving miscommunication, as simple dialogue systems are often incapable of doing so [25].

365 In general, however, the question arises as to whether this distribution of roles makes sense. Sauppé and Mutlu [21]  
366 work shows that users want robots to help them correct failures.  
367

368 *For which kind of failure is communication & fixing appropriate?* In addition to the responsibility for fixing the  
369 failure, the question arises: For which types of failure and in which situations are communication and fixing of failures  
370 beneficial? Especially in the case of repeated failures known to users, excessive communication by means of explanations  
371 could be counterproductive and lead to frustration among users. Moreover, not all failures can be fixed by users. As  
372 Honig and Oron-Gilad [12] highlights, whether a user can fix a failure depends on their skills, if the failure is perceived  
373 as such and if the user is motivated to fix it. A survey asking employees about the use of AI in their companies shows  
374 that the development of AI skills in employees is a challenge for companies to successfully use AI technologies [1].  
375 The failure can be successfully resolved only when a user has the knowledge and skills about the robot and how it  
376 works. Even if there is a step-by-step solution to fix the failure, the user must understand and be able to implement the  
377 individual steps.  
378

381 *How should the failure be communicated?* A failure can only be corrected if it is recognised. Therefore, the question  
382 arises: How should a failure be communicated so that users perceive it and are motivated to fix it? Research provides  
383 several possible design approaches and studies. Cha et al. [3], for example, pointed out that sound can transport the need  
384 for help and that light can help to communicate the urgency. The challenge in designing failure communication is that  
385 it must be appropriate to the application and situation. For example, a loud, continuous sound from a malfunctioning  
386 service robot in a meeting would be inappropriate, making it very difficult to continue the meeting. However, such  
387 a loud sound could be helpful in a production hall, as there is often a loud noise level, and the information would  
388 otherwise not be perceived. However, alternative forms of communication, e.g. notifying workers individually or using  
389 light signals, could also be helpful communication tools.  
390

393 *What about the robot-centered perspective?* Based upon the elaboration so far, it can be assumed that we perceive the  
394 human-centered perspective to be a vital aspect of human-robot interaction. Not taking the human in the centre ought  
395 to result in technical devices with a valuable purpose and with a user-unfriendly interface making their usefulness  
396 questionable. However, the technical system's or, in particular, the robot's perspective must also be taken into account  
397 when preventing, communicating, and resolving failures. Using the example of Hald et al. [10] and our proposed  
398 study idea: How shall the robot know of the occurrence of a failure? What indicates the failure? How can the robot  
399 differentiate between successfully having sorted the bottles and not completing its task? Given it knows of the failure,  
400 how shall it know it is a computer vision error due to bad lighting conditions and not a misplaced bottle? Given it  
401 knows the cause of the failure, how shall it know that turning on the light can fix the failure and not move its own  
402 position to move its shadow of the bottles?  
403

404 Only the symbiotically linking of both perspectives can create a trustful interaction between humans and robots.  
405

## 408 REFERENCES

- 410 [1] Elisabeth André, Wilhelm Bauer, Martin Braun, Chi Tai Dang, Matthias Peissner, and Katharina Weitz. 2021. *KI-Kompetenzentwicklung bei Sach und Produktionsarbeit*. Plattform Lernende Systeme. 1–34 pages.
- 411 [2] Daniel J Brooks. 2017. *A human-centric approach to autonomous robot failures*. Ph.D. Dissertation.
- 412 [3] Elizabeth Cha, Maja Matarić, and Terrence Fong. 2016. Nonverbal signaling for non-humanoid robots during human-robot collaboration. In *2016 11th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*. 601–602. <https://doi.org/10.1109/HRI.2016.7451876> ISSN: 2167-2148.
- 413 [4] Devleena Das, Siddhartha Banerjee, and Sonia Chernova. 2021. Explainable AI for robot failures: Generating explanations that improve user assistance in fault recovery. In *Proceedings of HRI '21: ACM/IEEE International Conference on Human-Robot Interaction, Boulder, CO, USA, March*

- 417        8-11, 2021, Cindy L. Bethel, Ana Paiva, Elizabeth Broadbent, David Feil-Seifer, and Daniel Szafrir (Eds.). ACM, New York, NY, USA, 351–360.  
418        <https://doi.org/10.1145/3434073.3444657>
- 419        [5] Ewart J de Visser, Marieke MM Peeters, Malte F Jung, Spencer Kohn, Tyler H Shaw, Richard Pak, and Mark A Neerincx. 2020. Towards a theory of  
420        longitudinal trust calibration in human–robot teams. *International Journal of Social Robotics* 12, 2 (2020), 459–478. <https://doi.org/10.1007/s12369-019-00596-x>
- 421        [6] Sahba El-Shawa, Noah Kraemer, Sara Sheikholeslami, Ross Mead, and Elizabeth A. Croft. 2017. “Is this the real life? Is this just fantasy?”: Human  
422        proxemic preferences for recognizing robot gestures in physical reality and virtual reality. In *2017 IEEE/RSJ International Conference on Intelligent  
423        Robots and Systems (IROS)*. 341–348. <https://doi.org/10.1109/IROS.2017.8202178> ISSN: 2153-0866.
- 424        [7] Connor Esterwood and Lionel P. Robert Jr. 2023. Three Strikes and you are out!: The impacts of multiple human–robot trust violations and repairs  
425        on robot trustworthiness. *Computers in Human Behavior* 142 (May 2023), 107658. <https://doi.org/10.1016/j.chb.2023.107658>
- 426        [8] Romi Gideoni, Shabee Honig, and Tal Oron-Gilad. 2022. Is It Personal? The Impact of Personally Relevant Robotic Failures (PeRFs) on Humans’  
427        Trust, Likeability, and Willingness to Use the Robot. *International Journal of Social Robotics* (Sept. 2022). <https://doi.org/10.1007/s12369-022-00912-y>
- 428        [9] David Gunning and David Aha. 2019. DARPA’s explainable artificial intelligence (XAI) program. *AI Magazine* 40, 2 (2019), 44–58. <https://doi.org/10.1609/aimag.v40i2.2850>
- 429        [10] Kasper Hald, Katharina Weitz, Elisabeth André, and Matthias Rehm. 2021. “An Error Occurred!” - Trust Repair With Virtual Robot Using Levels of  
430        Mistake Explanation. In *Proceedings of the 9th International Conference on Human-Agent Interaction (HAI ’21)*. Association for Computing Machinery,  
431        New York, NY, USA, 218–226. <https://doi.org/10.1145/3472307.3484170>
- 432        [11] Peter A Hancock, Deborah R Billings, Kristin E Schaefer, Jessie YC Chen, Ewart J De Visser, and Raja Parasuraman. 2011. A meta-analysis of factors  
433        affecting trust in human–robot interaction. *Human factors* 53, 5 (2011), 517–527.
- 434        [12] Shabee Honig and Tal Oron-Gilad. 2018. Understanding and resolving failures in human–robot interaction: Literature review and model development.  
435        *Frontiers in Psychology* 9 (2018), 861. <https://doi.org/10.3389/fpsyg.2018.00861>
- 436        [13] John D Lee and Katrina A See. 2004. Trust in automation: Designing for appropriate reliance. *Human factors* 46, 1 (2004), 50–80.
- 437        [14] Martina Mara, Jan-Philipp Stein, Marc Erich Latoschik, Birgit Lugrin, Constanze Schreiner, Rafael Hostettler, and Markus Appel. 2021. User  
438        Responses to a Humanoid Robot Observed in Real Life, Virtual Reality, 3D and 2D. *Frontiers in Psychology* 12 (2021). <https://www.frontiersin.org/articles/10.3389/fpsyg.2021.633178>
- 439        [15] Daxton Mitchell, HeeSun Choi, and Justin M. Haney. 2020. Safety Perception and Behaviors during Human-Robot Interaction in Virtual Environments.  
440        *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 64, 1 (Dec. 2020), 2087–2091. <https://doi.org/10.1177/1071181320641506>
- 441        [16] Cecilia G. Morales, Elizabeth J. Carter, Xiang Zhi Tan, and Aaron Steinfeld. 2019. Interaction Needs and Opportunities for Failing Robots. In  
442        *Proceedings of the 2019 on Designing Interactive Systems Conference (DIS ’19)*. Association for Computing Machinery, New York, NY, USA, 659–670.  
443        <https://doi.org/10.1145/3322276.3322345>
- 444        [17] James Reason. 2000. Human error: models and management. *Bmj* 320, 7237 (2000), 768–770.
- 445        [18] Mark O Riedl. 2019. Human-centered artificial intelligence and machine learning. *Human Behavior and Emerging Technologies* 1, 1 (2019), 33–36.  
446        <https://doi.org/10.1002/hbe2.117>
- 447        [19] Alessandra Rossi, Kerstin Dautenhahn, Kheng Lee Koay, and Michael L. Walters. 2017. How the Timing and Magnitude of Robot Errors Influence  
448        Peoples’ Trust of Robots in an Emergency Scenario. In *Social Robotics (Lecture Notes in Computer Science)*, Abderrahmane Kheddar, Eiichi Yoshida,  
449        Shuzhi Sam Ge, Kenji Suzuki, John-John Cabibihan, Friederike Eyssel, and Hongsheng He (Eds.). Springer International Publishing, Cham, 42–52.  
450        [https://doi.org/10.1007/978-3-319-70022-9\\_5](https://doi.org/10.1007/978-3-319-70022-9_5)
- 451        [20] Maha Salem, Gabriella Lakatos, Farshid Amirabdollahian, and Kerstin Dautenhahn. 2015. Would you trust a (faulty) robot? Effects of error, task  
452        type and personality on human–robot cooperation and trust. In *Proceedings of the Tenth Annual ACM/IEEE International Conference on Human-Robot  
453        Interaction, HRI 2015, Portland, OR, USA, March 2–5, 2015*, Julie A. Adams, William D. Smart, Bilge Mutlu, and Leila Takayama (Eds.). ACM, 141–148.  
454        <https://doi.org/10.1145/2696454.2696497>
- 455        [21] Allison Sauppé and Bilge Mutlu. 2015. The Social Impact of a Robot Co-Worker in Industrial Settings. In *Proceedings of the 33rd Annual ACM  
456        Conference on Human Factors in Computing Systems (CHI ’15)*. Association for Computing Machinery, New York, NY, USA, 3613–3622. <https://doi.org/10.1145/2702123.2702181>
- 457        [22] Raymond Ka-Man Sheh. 2017. “Why Did You Do That?” Explainable Intelligent Robots. In *The Workshops of the The Thirty-First AAAI Conference on  
458        Artificial Intelligence, Saturday, February 4–9, 2017, San Francisco, California, USA (AAAI Technical Report, Vol. WS-17)*. AAAI Press.
- 459        [23] Youssef Shibani, Julia Diemer, Simone Brandl, Rebecca Zack, Andreas Mühlberger, and Stefan Wüst. 2016. Trier Social Stress Test in vivo and  
460        in virtual reality: Dissociation of response domains. *International Journal of Psychophysiology* 110 (Dec. 2016), 47–55. <https://doi.org/10.1016/j.ijpsycho.2016.10.008>
- 461        [24] Suzanne Tolmeijer, Astrid Weiss, Marc Hanheide, Felix Lindner, Thomas M. Powers, Clare Dixon, and Myrthe L. Tielman. 2020. Taxonomy of  
462        Trust-Relevant Failures and Mitigation Strategies. In *Proceedings of the 2020 ACM/IEEE International Conference on Human-Robot Interaction (HRI  
463        ’20)*. Association for Computing Machinery, New York, NY, USA, 3–12. <https://doi.org/10.1145/3319502.3374793>
- 464        [25] Katharina Weitz, Lindsey Vanderlyn, Thang Vu Ngoc, and Elisabeth André. 2021. It’s our fault!: Insights Into Users’ Understanding and  
465        Interaction With an Explanatory Collaborative Dialog System. In *Proceedings of the 25th Conference on Computational Natural Language Learning,  
466        CoNLL 2021, Online, November 10–11, 2021*, Arianna Bisazza and Omri Abend (Eds.). Association for Computational Linguistics, 1–16. <https://doi.org/10.18653/v1/2021.conll-1.1>

- 469 [26] A. B. Ünal, R. Pals, L. Steg, F. W. Siero, and K. I. van der Zee. 2022. Is virtual reality a valid tool for restorative environments research? *Urban  
470 Forestry & Urban Greening* 74 (Aug. 2022), 127673. <https://doi.org/10.1016/j.ufug.2022.127673>
- 471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520