

# 1 Analyzing Trust in AI: 2 Factors in Safety-critical and non-safety-critical Situations

3  
4 *AI systems are complex and difficult to understand, which makes it essential to*  
5 *establish a basis of trust between the user and the system to enable comfortable*  
6 *interactions. AI systems present both opportunities and risks, hence developers*  
7 *should endeavor to minimize risks as much as possible. Users, particularly in the*  
8 *realm of safety-critical systems, require sufficient trust in any new technology to*  
9 *consider its usage. Uncertainty and missing understanding about the field of AI*  
10 *and algorithms in general hinders and stops its usage. Therefore, transparency*  
11 *and comprehensibility should be goals in the development of AI systems. These*  
12 *factors can be divided into three groups: Technical factors, social factors, and*  
13 *evaluation factors. Each of these groups consists of five factors that can either*  
14 *enhance or diminish trust in an AI system. Based on these factors, a reference*  
15 *model is created to classify these three groups of factors within the overall context*  
16 *of trust formation. A study was conducted with 258 participants using a*  
17 *questionnaire to examine and subsequently analyze and evaluate the previously*  
18 *formulated reference model. An attempt is made to hierarchize the individual*  
19 *factors within the three groups and to identify any differences related to safety-*  
20 *critical and non-safety-critical systems.*

21  
22 **Keywords:** *artificial intelligence, trustworthy ai, responsible ai*

## 23 24 25 Introduction

26  
27 Machine learning, particularly deep neural networks, are often challenging to  
28 comprehend. A known input value is transformed into a potential output value  
29 predetermined by the network developer's design decisions which control the  
30 complex computations within the network. The structure of the network, with its  
31 various layers, is specified by the developer. However, the precise calculations and  
32 intermediate values that an input value undergoes on its path to becoming an output  
33 value are unknown. Even developers of highly complex neural networks typically  
34 only know what type of operations are performed in each respective layer. These  
35 networks are referred to as black boxes, as the exact behaviors of the neural network  
36 cannot be determined or described. Consequently, developers of such neural  
37 networks must focus particularly on input data and its structure. Input data requires  
38 special attention, as the exact processing within the network is unknown. For  
39 instance, it is challenging to identify sample biases within the input data that may  
40 render the results of the neural network unusable. Such biases could occur due to  
41 unnoticed characteristics within the data.

42 Users of AI systems thus desire transparency and a certain degree of  
43 comprehensibility, especially when employing machine learning algorithms in  
44 safety-critical application domains, such as medical technology or autonomous  
45 vehicles. Physicians receiving a diagnosis from an algorithm can better assess the  
46 accuracy of the diagnosis if they have background information about the algorithm's  
47 decision-making process. Furthermore, a more thorough understanding of a black

1 box can aid in error corrections or the debugging process. People typically fear the  
2 unknown, and similar sentiments apply to trusting something opaque (Carleton,  
3 2016). Explainable artificial intelligence approaches aim to counter this issue, as  
4 more and more systems are equipped with AI functionalities, and trust in AI systems  
5 should be strengthened. Trust plays a significant role in the utilization of  
6 mechanisms and devices whose functionalities cannot be fully understood by its  
7 users. This trust can be gained through test results or personal experiences. Whether  
8 it concerns the use of a technical device or trust in another individual, it is more  
9 about an inner conviction, typically based on positive experiences. Thus, the concept  
10 of trust can always be associated with a dynamic relationship, which can be  
11 established with both humans and objects.

12 According to Luhmann, trust can be used to reduce complexity in modern  
13 societies. Uncertainties and dangers can be weakened. Applied to AI systems, this  
14 means that a certain level of trust is required to interact meaningfully with a complex  
15 system. The lack of information about the system must be compensated for and  
16 bridged with the help of trust (Luhmann, 1968). In interpersonal relationships, trust  
17 plays a crucial role. This trust must be built over an extended period and requires  
18 constant attention. Trust in another person can quickly be compromised. When  
19 using algorithms, it is expected that their use should occur directly and without a  
20 prolonged period of trial and adjustments. Experiences are typically relayed second-  
21 hand, through consultations, or from specialized personnel. Trust can be challenging  
22 to establish under such conditions. Personal experiences can only be gained through  
23 active use, which is not an effective way to build trust, particularly with novel and  
24 safety-critical systems. Even with non-critical systems, such as using a chatbot,  
25 users often lack the necessary information to understand how the system generates  
26 a response.

27  
28

## 29 **Literature Review**

30

31 In terms of trust in AI systems, there are already numerous research papers and  
32 scientific articles. The European Commission sets an important benchmark for the  
33 development of AI systems in Europe with their "Ethics Guidelines for Trustworthy  
34 AI" (European Commission, 2019). Accordingly, trustworthy AI is characterized by  
35 seven components. The expert group identifies human agency, robustness, privacy,  
36 transparency, ethical behavior, sustainability, and accountability as the most  
37 important components in the development of trustworthy AI systems. The expert  
38 group of the European Commission focuses on the conception of AI systems and  
39 describes factors that form a basis for trustworthy AI systems. Without these  
40 fundamental factors, trust in AI systems is not possible. Ideally, all three factors  
41 should apply to an AI system for it to be considered trustworthy. But other factors  
42 such as explainability, as well as security, controllability, reliability, and experience,  
43 should not be underestimated (Jacovi, 2021). These factors are equally important  
44 when it comes to creating genuine trust in an AI system. The more positive factors  
45 associated with a system, the easier it is to build trust in that system.

1 In view of this multitude of factors, this work divides these factors into two  
2 groups. On the one hand, there are technical factors, and on the other hand, there are  
3 social factors. Technical factors are not directly transferable to a trust relationship  
4 between two individuals, whereas this is the case for social factors. Hoff et al. write  
5 that relationships between humans and machines as well as between two humans  
6 are different but fundamentally based on the same principles (Hoff, 2015). Taking a  
7 closer look at the factors of an interpersonal relationship, factors such as reliability  
8 or traceability emerge, which can also be transferred to trust in an AI system. The  
9 foundation of trust outlined by the expert group of the European Commission can  
10 thus be compared to a first impression of another person. Trust itself can only be  
11 built and strengthened as the relationship progresses, offering much scope for  
12 considering the additional factors for a more detailed analysis. Factors such as  
13 transparency and explainability play a prominent role here, while factors such as  
14 reliability or security are also important but may be weighed differently by various  
15 users.

16 Parasuraman et al. mention an interesting observation here. It was found that  
17 people continue to use a system after an error has occurred, even without  
18 understanding the error (Parasuraman, 1997). Such behavior would be atypical in  
19 an interpersonal relationship and is likely more common in the use of computer  
20 systems. Without an explanation of why something did not work, few people would  
21 probably trust another person again. Parasuraman et al. suggest a dependency  
22 between humans and machines here, as, for example, an activity may not be possible  
23 without the assistance of the machine. If such a dependency exists, people are also  
24 willing to continue working with a system even in the case of serious errors.  
25 Consequently, trust can be replaced by dependencies. However, according to  
26 Ferronato et al., a lack of trust and forced interaction with a system, such as in a  
27 work context, lead to a reduction in effectiveness and safety (Ferronato, 2020). In  
28 the context of interpersonal relationships, such dependence could be compared to  
29 coercion and thus does not reflect a desirable state. Trust should therefore always be  
30 an important goal in the design of AI systems, while efforts should be made to avoid  
31 usage due to dependency.

32 Ferronato et al. also address another interesting aspect. Trust itself can be  
33 further subdivided into "Primary Trust" and "Secondary Trust" (Ferronato, 2020).  
34 "Primary Trust" refers to the initial contact with something or someone, the first  
35 impression. "Secondary Trust", on the other hand, entails the trust that develops over  
36 time. It has been observed that the first impression has a very strong influence on  
37 "Secondary Trust". Furthermore, personal predisposition, which varies from person  
38 to person, plays a crucial role in the formation of "Primary Trust". Factors such as  
39 prior experiences, experiences of third parties, or the reputation of the system are  
40 important here (Ferronato, 2020). Trust in an AI system is therefore significantly  
41 based on the first impression, which is also strongly influenced by the fundamental  
42 attitude towards AI systems in general. Thus, one could speak of a "three-step  
43 process" (fundamental attitude, first impression, usage) of trust, as usage is highly  
44 unlikely to occur in safety-critical domains with a negative fundamental attitude or  
45 a negative first impression.

1 Additional characteristics of trust were discovered by Alarcon et al. They report  
2 on the dangers of "over-trusted" and "under-trusted" systems (Alarcon, 2020). If a  
3 system is trusted too much, the capabilities and functionalities of a system are  
4 overestimated. This leads to the danger that users interact with a system too casually  
5 and, for example, do not perform quality checks on the results of an AI system. It is  
6 assumed here that the system will provide correct results anyway. Conversely, if a  
7 system is trusted too little, functionalities of a system are not utilized because it is  
8 assumed from the outset that the system will not provide useful results (Alarcon,  
9 2020). Both scenarios are suboptimal as they result in the underutilization of a  
10 system's full potential. Establishing trust in a system is crucial, but it's equally  
11 important not to overestimate its capabilities. Consequently, developers of AI  
12 systems should incorporate trust considerations into the development process. Trust  
13 is elusive and demands a considerable degree of attention and deliberation.

14 But how does trust arise and how can a trusting relationship be described?  
15 Hancock et al. write that the principle of trust is always based on three elements.  
16 These include two agents and a communication channel. One of the agents is the  
17 sender, and the other is the receiver (Hancock, 2011). The sender is the agent who  
18 transmits trust, that is, who trusts the other agent. The receiver is trusted by the  
19 sender; thus, the receiver receives trust through the communication channel. In the  
20 context of AI systems, the communication channel can be referred to as the  
21 operational environment. This could be one's own home computer for a software  
22 application or one's own car for a hardware application. Furthermore, Hancock et  
23 al. write that trust is a dynamic variable (dimension) that becomes increasingly  
24 important with increasing autonomy (Hancock, 2011). The concept of a humanoid  
25 robot, as depicted in science fiction films, would represent one possible stage of this  
26 dimension, as trust in such a robot would be nearly equivalent to that in another  
27 human. However, all intermediate stages and developments of current AI systems  
28 need to be considered in a more nuanced manner. They are gradually approaching a  
29 classical trust relationship between two humans as AI systems become more  
30 advanced and autonomous in decision-making.

31 In reference to the previously mentioned image of two agents, Devitt writes  
32 that the sender always makes themselves vulnerable, as they trust that the receiver  
33 will not exploit them (Devitt, 2018). This fundamental fear applies only partially to  
34 AI systems, as it involves the "user - application" situation. The user tends to trust  
35 that the system will perform a task correctly rather than trusting that the system will  
36 perform any task at all. Therefore, in the context of AI systems, special attention  
37 must be paid to the level of functional trust. Functional trust describes the trust in a  
38 system's ability to successfully execute a specific task. Consequently, the focus in  
39 further examination is placed more on the topic of functional trust, while  
40 interpersonal factors are considered only marginally. Hancock et al. address an  
41 important point here. They highlight the role of "false trust" in the development of  
42 trustworthy AI systems. False trust involves the concept of deception. Through  
43 deception, a system may appear trustworthy or successful in its actions when it is  
44 not (Hancock, 2011). Such a system could, for example, secretly collect data, and  
45 the expected task is merely a secondary activity to appear trustworthy. If this fact  
46 were to be revealed, no one would knowingly interact with the system anymore

1 because without "false trust", it would not enjoy genuine trust. However, there is  
2 also a danger here if system functions offer significant advantages or substantial  
3 workflow simplification. For example, it might be possible for AI systems to be  
4 used, especially in non-critical areas, without being trusted to a high degree. Devitt  
5 mentions companies as an example of the relevance of this consideration, which do  
6 not prioritize data protection and share personal data to increase their own revenue  
7 (Devit, 2018). In this regard, there are interesting questions for further investigation.  
8 How significant must the functional advantage be for a person to interact with a non-  
9 trustworthy AI system? Functional trust plays a crucial role in AI systems and can  
10 be enhanced through transparency and explainability of algorithms. This is  
11 important, because the potential underuse of AI technology could be one logical  
12 consequence of the fear of misuse (Floridi, 2018).

13 To better understand the concept of deception in relation to trust, it is important  
14 to consider the opposite of trust. After a brief consideration, one may think of the  
15 term mistrust. But is mistrust truly the opposite of trust? Smithson writes that a more  
16 nuanced examination is required to answer this question (Smithson, 2018).  
17 According to Smithson, mistrust and the absence of trust are not synonymous. In  
18 the case of mistrust, clear doubts about trustworthiness can be identified based on  
19 information. In the case of the absence of trust, there is not enough information  
20 available to make a definitive decision about whether the system is trustworthy or  
21 not. However, the absence of trust is clearly the scenario that should primarily be  
22 considered when developing new AI systems. Mistrust is too strong a term when it  
23 is a matter of the absence of trust. Deception and mistrust should therefore be  
24 identified as extreme values.

25 In considering functional trust, the focus is on establishing trust in the  
26 functionalities of an AI system. Individuals lacking an understanding of a system's  
27 functions must be able to place their trust in the system's functionalities to utilize it  
28 efficiently. The process of building trust is therefore essential for the utilization of a  
29 system. As this process varies from system to system, preconceptions should be  
30 made during the development process.

31

### 32 *The Human Black Box*

33

34 The preceding examination of trust between humans and AI systems raises  
35 some questions regarding the human decision-making process. When one person  
36 trusts another person, they are also trusting a black box. When a person interacts  
37 with an AI system, they desire transparency and an understanding of its behavior.  
38 This paradox has also been addressed and analyzed by Bonezzi et al. They discuss  
39 the illusion that the decision-making process of a human is more transparent than  
40 that of an AI system (Bonezzi, 2022). A person is more likely to trust another person  
41 because they see themselves as a template and can more easily transfer their own  
42 decision-making process onto other people. An AI algorithm is foreign, and one's  
43 own decision-making process cannot be readily transferred onto an algorithm.

44 Many regulations are being introduced to demand transparency of algorithms,  
45 such as those by the European Union, to provide people with a better understanding  
46 of AI systems (Goodman, 2016). The intentions behind these mechanisms are clear.

1 However, daily, humans interact with a multitude of black boxes that are not  
2 required to openly disclose why they behave as they do. Examples of this include  
3 doctors (Bonezzi, 2022). Diagnoses are made without further explanation to a  
4 patient. It is trusted that a physician makes a correct diagnosis, but misdiagnoses  
5 occur. However, there are no initiatives to make human decisions explainable and  
6 transparent. The decisions of a taxi driver are as much a black box as those of an  
7 autonomous vehicle. Trust is established through laws, rules, and the transferability  
8 of one's own decision-making process to the taxi driver, even though no trust  
9 relationship has been built with an individual. This concept of trustworthiness is an  
10 important dimension when dealing with an uncertain future (Kusche, 2023). Kusche  
11 describes it as bridging the distinction between decision-makers and those affected.  
12 An AI system does not make conscious decisions and is therefore a foreign entity to  
13 many people when considering trusting something rather than someone.

14 Bonezzi et al. have found that people tend to empathize with other people,  
15 while this is not the case with AI systems. When people are closer to each other, this  
16 behavior occurs more frequently, and trust can be built more easily (Bonezzi, 2022).  
17 People belonging to a cohesive group are more likely to trust each other than those  
18 without close contact. However, in both cases, trust is placed in a black box, whether  
19 it involves friends or strangers. The human being uses themselves as a template to  
20 trust another human. The question that arises is: If an algorithm, which is a white  
21 box, could be used as a template to illustrate the decision-making process of another  
22 algorithm, which is a black box, would a person trust this algorithm?

23 To answer this question, three types of factors are distinguished: technical,  
24 social, and evaluation factors. Technical factors and social factors are no new  
25 concept, other works have already established similar kinds of factors. Ikkatai et al.  
26 have created an octagon metric with eight important factors for formulating trust  
27 (Ikkatai, 2022). Kaplan et al. and Bach et al. also worked on the identification  
28 process for factors in which trust can be measured (Kaplan, 2021; Bach, 2022). The  
29 evaluation factors are an additional set of factors that this work introduces to better  
30 understand outside influences that can determine the process of forming trust.

31

32

### 33 **Methodology**

34

35 The following section will focus on the factors influencing trust. These factors  
36 were decided upon reviewing other works in this field. Additionally, a set of nine  
37 test runs were conducted before the actual study took place to validate the approach.  
38 This led to the sub factors in the next part for generating a better understanding of  
39 main factors for participants of the study.

40 Regarding the technical factors, the AI system itself is considered. This includes  
41 stability, applicability, behavior, comprehensibility, and accuracy. These factors  
42 consist of several sub-factors. Technical factors are distinguished from social and  
43 evaluation factors in that they cannot be easily transferred to a trust relationship  
44 between two people. For instance, a person's behavior might be attributed to their  
45 political or religious beliefs. Since technical factors are intended to form the basis

1 for social factors, it was decided to include data protection and ethical correctness  
2 in the category of technical factors.

- 3
- 4 • Stability (Robustness/Responsiveness)
- 5 • Applicability (Usability/Difficulty)
- 6 • Behavior (Privacy/Ethics)
- 7 • Comprehensibility (Transparency/Explainability)
- 8 • Accuracy (Results/Reliability)
- 9

10 The second group of factors comprises social factors. These are characterized  
11 by their resemblance to elements found in human trust relationships. They are easier  
12 to transfer than technical factors and consider the AI system as an entity worthy of  
13 trust. The technical implementation is consequently functional and instills basic trust  
14 in the form of the factors already presented. Social factors include traceability,  
15 reliability, accountability, overall competence, and transferability. Transferability  
16 refers to the ability to transfer one's own values onto the AI system.

- 17
- 18 • Traceability
- 19 • Reliability
- 20 • Accountability
- 21 • Competence
- 22 • Transferability
- 23

24 The third group of factors comprises evaluation factors. These factors  
25 encompass the overarching theme of AI systems and societal perception. They  
26 examine reputation as well as experiential values. This specific group of factors aims  
27 to illustrate the influence of these values on trust in AI systems. It considers the  
28 media and popular culture reputation as well as the reputation of the developer or  
29 manufacturer. Furthermore, it investigates whether reviews, experiences with other  
30 AI systems, or the experiences of other individuals impact trust in AI systems.

- 31
- 32 • Reputation (Media/Popular culture)
- 33 • Reputation (Developer/Manufacturer)
- 34 • Experiences (Other AI systems)
- 35 • Experiences (Other people)
- 36 • Testing reports (Magazines/Consultations)
- 37

### 38 *Introducing the Trust Pyramid as a Hierarchy*

39

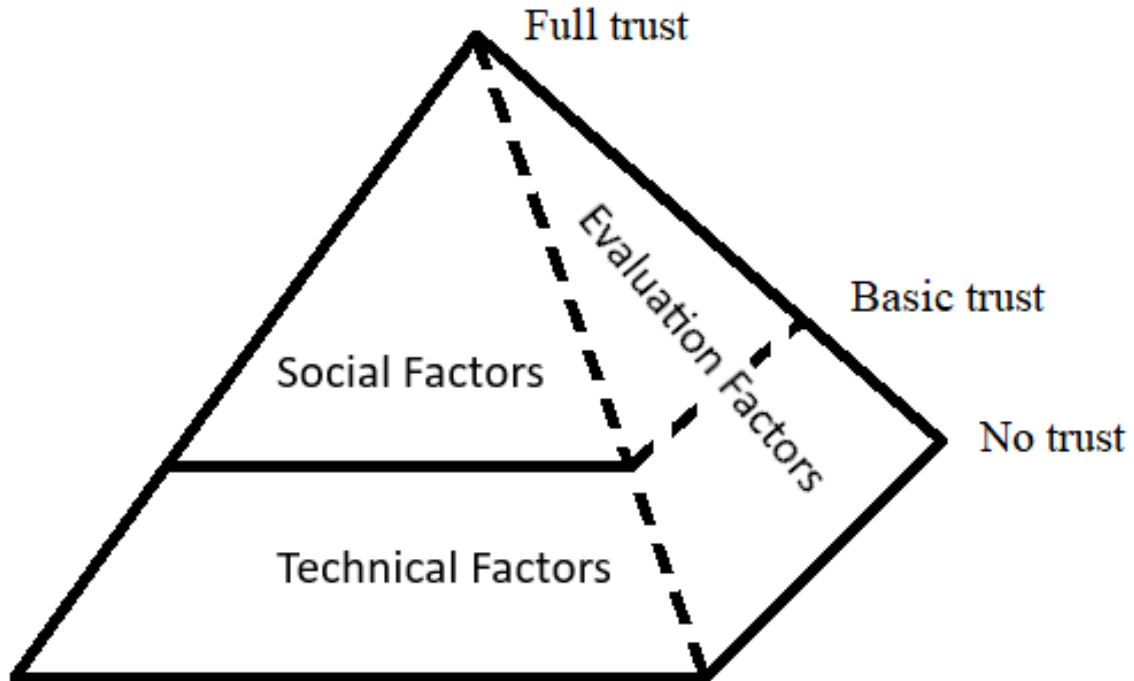
40 A study was conducted to evaluate the aforementioned factors and find answers  
41 to the following two questions:

- 42 1. What factors hold the highest priority when it comes to fostering trust in an  
43 AI system?
- 44 2. Does the context in which an AI system makes its decision play a crucial  
45 role?

1  
2  
3  
4  
5  
6  
7  
8  
9

The first question concerns the factors or groups of factors that have already been identified, while the second question relates to the context of use, with a focus on security aspects. Therefore, when both questions are combined, two scenarios are emphasized: one involving an AI system in a non-safety-critical context and the other involving an AI system in a safety-critical context. Figure 1 illustrates the trust pyramid that has been developed, which relates the trust factors to each other and attempts to classify the factors in terms of overall trust in any given system.

10 **Figure 1.** *Trust Pyramid - Based on own representation*



11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

The trust pyramid always refers to a single separately considered AI system from the perspective of a singular individual. The basis on which trust is forming is different for everyone. Because of this the process of trusting a system or a person is not generalized by this approach. This work assumes a personal fundamental attitude of each human being. This varies from person to person and is based on experiences, one's own knowledge base, and personal worldview. If this fundamental attitude is not yet developed or if someone has no idea what an AI system is, it's a different process then for someone who has broad experiences with different AI systems. At the base of the pyramid, someone has no trust in a specific AI system. This is generally the case with any AI system, as there has been no engagement with it yet.

The bottom layer of the actual pyramid represents the technical factors. Here, a foundational trust is established based on stability, behavior, comprehensibility, accuracy of results, and usability. Without these factors, social factors don't play any role in generating trust, as trust in an AI system is not possible without trust in the technology. A broken system that does not work correctly or does not produce



1 results will have a hard time generating trust. The same is the case for any given  
2 system that can't be understood or behaves in a way that is not making sense to its  
3 users. Therefore, this work assumes a dependency between these two groups of  
4 factors. The dashed line on the right side of the pyramid represents an overlap with  
5 the domain of the evaluation factors. It is assumed that even if a system is technically  
6 trustworthy, a significant portion of the foundational trust depends on trust in the  
7 developer or manufacturer. The reputation of a company or the political stance of a  
8 CEO can also have an impact. These factors are typically independent of the  
9 software or hardware components.

10 The social factors and evaluation factors fundamentally build upon the  
11 technical factors. Social factors such as transparency, reliability, accountability,  
12 competence, and transferability should be considered after establishing the technical  
13 framework. If, in addition to the technical factors, the social factors can also be  
14 fulfilled, strong trust in an AI system can be inferred. The left side of the pyramid  
15 composed of technical and social factors represents the maximum trust that can be  
16 achieved from an AI system. The right side of the pyramid is only attainable through  
17 additional trust in the evaluation factors. However, these factors may deviate even  
18 in a perfect system. For example, if negative news about AI systems is currently  
19 prevalent in the media, this perception may be transferred to other similar systems.  
20 Thus, complete trust is anchored at the absolute pinnacle of the trust pyramid and  
21 requires, in addition to a trustworthy AI system, a trustworthy developer/  
22 manufacturer and a generally positive perception of AI systems in the media or in  
23 the everyday environment.

#### 24 25 *Study Design*

26  
27 For the measurement within this study, a questionnaire is utilized. This  
28 questionnaire was developed based on the "OECD Guidelines on Measuring Trust"  
29 (OECD, 2017). As the conducted study is an online study, it mainly adhered to  
30 structural guidelines. Furthermore, no interviews were conducted. The study was  
31 held during a 46-day period with 258 active participants. The participants had to  
32 answer a questionnaire with questions related to different types of AI systems. Each  
33 question of the given scenarios was related to one of the factors.

34 For the completion of the questionnaire, a time span of 10 to 15 minutes was  
35 allocated, which was confirmed to be suitable through feedback from the respondents.  
36 The time frame was determined through a pre-test of the questionnaire. This pre-test  
37 also ensured that no technical difficulties would arise during the completion of the  
38 questionnaire. It was established beforehand what constitutes as "Valid cases" and  
39 what can be disregarded during data analysis. Additionally, this ensured that no  
40 necessary questions could be skipped and that dependencies within the  
41 questionnaire functioned correctly. Similarly, before the actual survey commenced,  
42 it was ensured that the questions were posed in a clear and comprehensible manner.

43 During the survey period, the questionnaire was completed by 294 individuals.  
44 However, 19 of these responses had to be excluded due to incomplete data.  
45 Consequently, 275 questionnaires were fully completed. Out of the 275 completed  
46 questionnaires, a further 17 were subsequently excluded due to erroneous responses

1 in the control questions. Thus, 258 datasets are relevant for the subsequent analysis  
2 of the results. 138 male, 118 female, 1 diverse and 1 participant without gender  
3 information were recorded. One objective of the study was to identify the factors  
4 that play the most significant role in fostering trust in an AI system.

## 5 6 7 **Results**

8  
9 A preliminary conclusion can be drawn regarding the initial question of this  
10 study. The objective of the question was to identify the factors that play the most  
11 significant role in fostering trust in an AI system. Concerning the technical factors,  
12 the factors of stability, applicability, and accuracy must be mentioned. These were  
13 identified in the study as the most important technical factors. For the social factors,  
14 also three factors can be identified: reliability, accountability, and competence. The  
15 factor of traceability can be somewhat related to the technical factor of  
16 comprehensibility. It has been highlighted that a complete understanding is not  
17 always necessary when using a system. The ability to review or question the  
18 decisions of an algorithm is sufficient for many participants. Furthermore, trust  
19 gained through other factors can be used to fill this "knowledge gap". The factor of  
20 transferability was challenging for many individuals to imagine. This was easier for  
21 questions regarding a chatbot compared to an autonomous vehicle. Specifically, this  
22 factor needs to be examined to determine when people desire human-likeness and  
23 when they do not.

24 In terms of the evaluation factors, the most significant factors that were  
25 identified were the reputation towards developers/manufacturers, the personal  
26 experience regarding AI systems, and professional test reports. While media and  
27 films/games can have an influence, objective, industry-specific articles, and  
28 opinions based on test reports have a higher significance when it comes to building  
29 trust. Similarly, the experiences of other individuals have a lesser impact compared  
30 to one's own experiences. The trust pyramid is intended to serve as the initial means  
31 of such a classification and to relate the three groups of factors to one another. In  
32 already mentioned works, all factors were considered equally. The study conducted  
33 as part of this work demonstrated that different systems require a varied weighting  
34 of individual factors within the factor groups. While stability and accuracy were  
35 more critical for autonomous vehicles, participants emphasized the behavior  
36 regarding privacy and ethics of the system concerning chatbots. Individuals are  
37 willing to disclose more data for safety-critical systems than for non-safety-critical  
38 systems. Conversely, safety-critical systems are expected to have a higher degree of  
39 accuracy and stability. These differences in expectations influence trust in a system  
40 and must therefore be given higher weight. A clear hierarchy of factors cannot  
41 therefore not be established without the context of the system in question.

42 The introduced groups of factors aim to provide a better description of AI  
43 systems and allow for the classification of an individual's trust in specific aspects of  
44 a system. The introduction of hierarchized factors adds to the approaches of Ikkatai  
45 et al., Kaplan et al. and Bach et al. by incorporating a factor of external impact,  
46 which can have a significant influence on trust formation (Ikkatai, 2022; Kaplan,

2021; Bach, 2022). The technical and social factors can also be found in those works but are complemented by the evaluation factors to categorize multiple factors. This helps facilitating the classification of interconnected aspects.

## Discussion

The second question of this study will now be addressed: "Does the context in which an AI system makes its decision play a crucial role?" An attempt will be made to find an answer to the issue just posed. If the context of an AI system does play a crucial role, then a clear sequence cannot be established across multiple systems. In the following, two different AI systems were compared. This involved distinguishing between a chatbot as a non-safety-critical system and an autonomous vehicle as a safety-critical system. To compare these two different scenarios, 17 questions from the previously presented factors were selected and posed again in relation to these two types of systems. Each question always refers to one of the three factor groups.

As mentioned in the previous section, the type of AI system likely plays a role in many factors. For example, differences in the complexity of tasks of an AI system or safety concerns can affect the creation of trust. An autonomous vehicle, for instance, affects the physical well-being of a user, while chatbots are currently deployed in less risky environments. Autonomous vehicles operate in an uncertain and unpredictable environment and therefore must solve very complex tasks accordingly. Furthermore, chatbots have been a more familiar technology for some time, especially after the introduction of ChatGPT 3 to the public. Autonomous vehicles, on the other hand, are a newer development and are not yet very prevalent in people's everyday lives. Media coverage and personal experiences are also possible criteria for differentiating the evaluation of these two systems.

To obtain a more precise picture of the situation, the same 258 participants were also asked about their trust in chatbots and autonomous vehicles. All questions were asked with a given scenario about the AI system. These questions were created with the OECD guidelines as a basis for measuring trust. In addition, the questions were also adjusted based on the test runs before the actual study took place.

1. For me to use the AI system, the system must be responsive.
2. For me to use the AI system, the system must be easy to use.
3. For me to use the AI system, the system must be limited to relevant data.
4. For me to use the AI system, the system must behave respectfully towards me.
5. For me to use the AI system, I must understand the system's procedures.
6. For me to use the AI system, the results must be precise and the predictions accurate.
7. For me to use the AI system, I must be able to inspect its decision-making process.
8. For me to use the AI system, the system must deliver reliable results.
9. For me to use the AI system, I must be able to assess the responsibilities.

- 1 10.For me to use the AI system, the system must competently complete
- 2 assigned tasks.
- 3 11.For me to use the AI system, the system must behave similarly to a human.
- 4 12.For me to use the AI system, it must be recommended to me by other people.
- 5 13.For me to use the AI system, the system must be positively reported in the
- 6 media.
- 7 14.For me to use the AI system, I must trust the developer/manufacturer.
- 8 15.For me to use the AI system, I must already trust other AI systems from the
- 9 same developer/manufacturer.
- 10 16.For me to use the AI system, many people must use the system.
- 11 17.For me to use the AI system, test reports must be positive.

12  
 13 Several T-tests were used to determine whether there are significant differences  
 14 between the two systems. For a more detailed representation of the differences, a  
 15 significant level of 0.001 was chosen. Table 1 provides an overview of the pairs and  
 16 whether the t-value is <0.001. Table 1 shows 17 pairs, each regarding one of the 15  
 17 factors listed under methodology. Two factors were split to get a better picture of the  
 18 two systems in comparison. These factors are behavior and experience. The factor's  
 19 behavior was split into privacy and ethical behavior, while the factor experience was  
 20 split into own experiences and experiences of other people.

21  
 22 **Table 1.** *T-test for comparing safety-critical and non-safety-critical AI systems*

Pairs	T	One-Sided p	Two-Sided p	Factor
Pair 1	-10.340	<.001	<.001	Stability
Pair 2	-4.059	<.001	<.001	Applicability
Pair 3	1.540	.062	.125	Behavior (Privacy)
Pair 4	1.895	.030	.059	Behavior (Ethics)
Pair 5	-8.766	<.001	<.001	Comprehensibility
Pair 6	-7.625	<.001	<.001	Accuracy
Pair 7	-9.579	<.001	<.001	Traceability
Pair 8	-6.387	<.001	<.001	Reliability
Pair 9	-10.092	<.001	<.001	Accountability
Pair 10	-7.433	<.001	<.001	Competence
Pair 11	-1.840	.033	.067	Transferability
Pair 12	-8.560	<.001	<.001	Reputation (Media)
Pair 13	-10.259	<.001	<.001	Reputation (Developer)
Pair 14	-9.864	<.001	<.001	Experience (AI systems)
Pair 15	-7.611	<.001	<.001	Experience (Other)
Pair 16	-9.679	<.001	<.001	Experience (Mine)
Pair 17	-10.900	<.001	<.001	Testing reports

23  
 24  
 25 From this table, it becomes evident that there are significant differences  
 26 between chatbots and autonomous vehicles in almost all factors. No differences  
 27 were observed in the factors "Behavior (Privacy)", "Behavior (Ethics)", and  
 28 "Transferability". For these three questions, it was important for the participating  
 29 individuals that both systems operate on the same principles. Significant differences  
 30 were detected for all other factors. Thus, it can be concluded that there is a varying

1 perception and evaluation of the aforementioned factors for most aspects.  
 2 Additionally, the negative T-value indicates that all significantly different values  
 3 regarding autonomous vehicles were rated more strongly. This confirms the  
 4 previously stated assumption that a clear hierarchy between trust factors cannot be  
 5 established independent of the AI system considered. Only the underlying AI system  
 6 provides the basis for considering the factors and enables individual prioritization  
 7 of each factor. Therefore, the trust pyramid can be used as a reference model to allow  
 8 for a rough categorization of the three groups of factors.

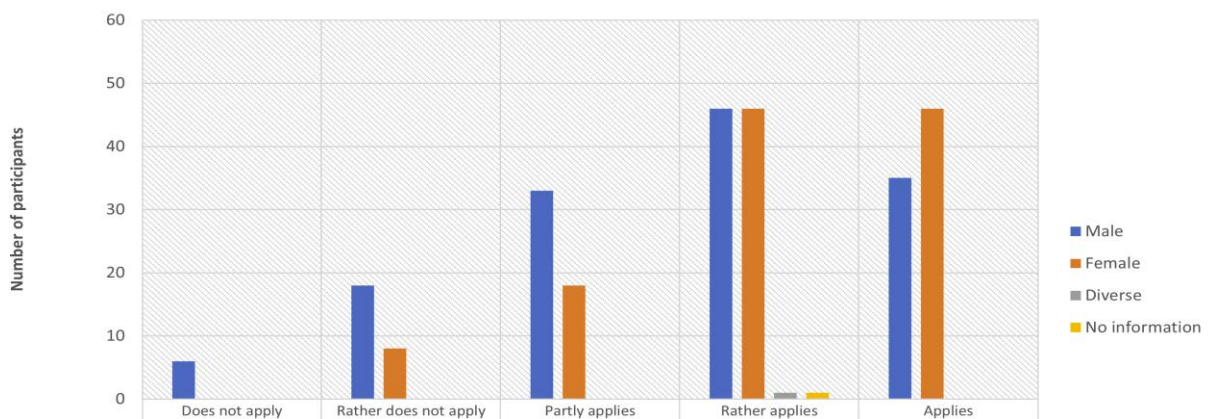
9 Based on the results of the T-test, the second question can be conclusively  
 10 answered with "Yes". The context in which an AI system makes its decisions plays  
 11 a crucial role. Trust in a safety-critical system is rated higher than in a non-safety-  
 12 critical system. Establishing a hierarchy of factors, as attempted at the end of the  
 13 previous section, is thus not possible and must be done individually when  
 14 considering a specific AI system. However, since the trust pyramid also includes  
 15 human considerations, a different order of factors can be assumed between two  
 16 different individuals for the same AI system. A specific examination of individual  
 17 factor groups concerning more than two AI systems could provide an approximation  
 18 of an order. This could identify types of AI systems that are like each other.

19 *Differences between female and male participants*

20 Due to the very similar number of female and male participants, differences  
 21 between these two groups will also be addressed. In general, it can be observed that  
 22 no significant differences were found between the two groups. To assess the  
 23 differences, a separate t-test was conducted for each group on an individual  
 24 question. Like the previous comparison between safety-critical and non-safety-  
 25 critical AI systems, a significance level of 0.001 was chosen.

26  
 27 **Figure 2.** *Diagram showing the single difference for technical factors*

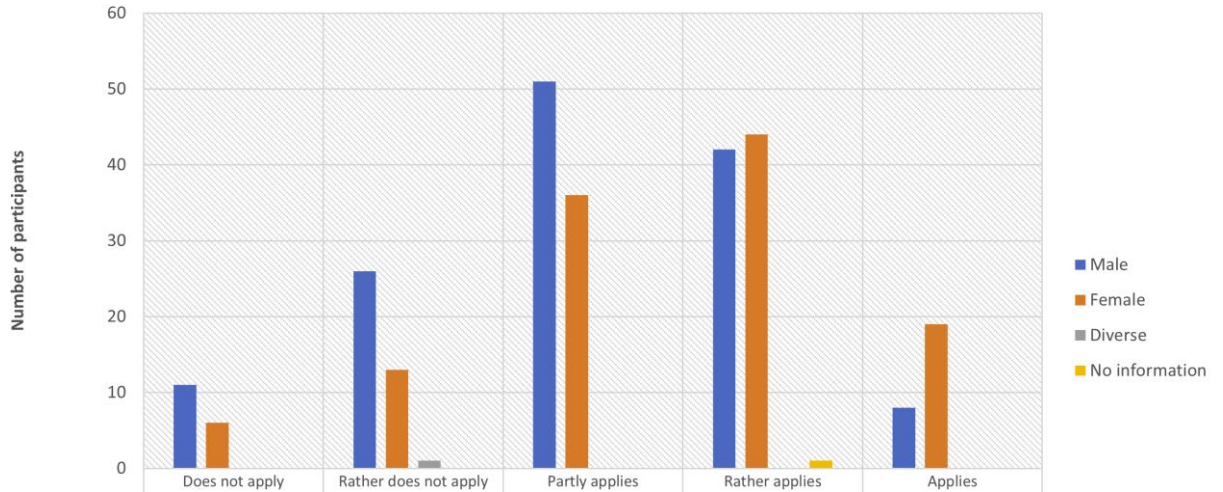
In order to use a chatbot, the system must behave respectfully towards me.



28  
 29  
 30 Figure 2 shows the only significant difference in the responses between men  
 31 and women regarding the technical factors. This question concerned the rapid  
 32 learnability of handling a system. While it was more important for female  
 33 participants in the study that the handling of an AI system is quickly learnable,  
 34 male participants indicated that it is not as important to them compared to females.

1 However, the other questions do not show significant differences between the two  
 2 comparison groups, which is why this difference warrants special mention.

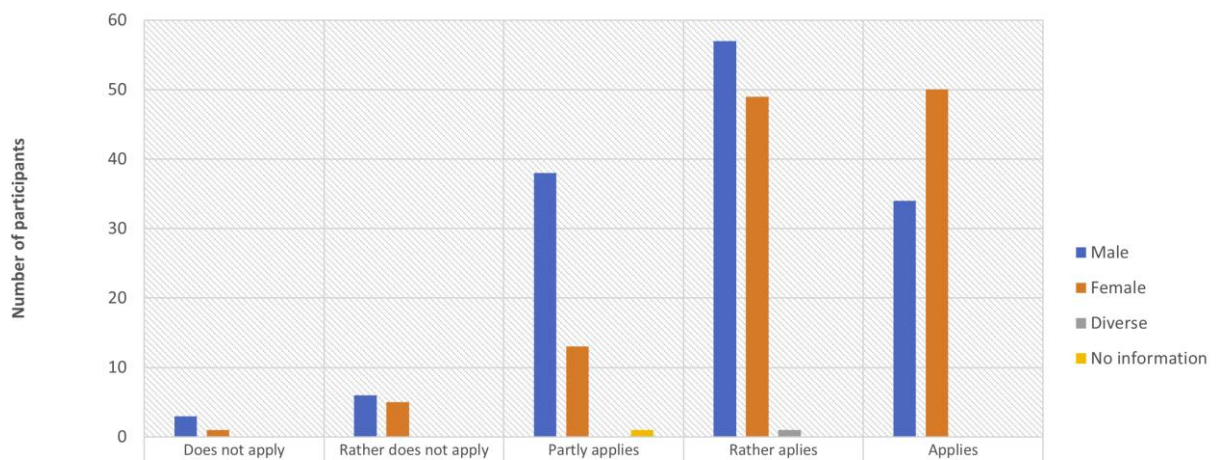
3  
 4 **Figure 3.** *Diagram showing the single difference for social factors*  
 I trust AI systems that have been recommended to me by other people.



5  
 6  
 7 Similarly, regarding the social factors, only one question exhibited a significant  
 8 difference between the genders. This difference is illustrated in figure 3. The  
 9 disparity between the two groups is evident in this question, which explored whether  
 10 trust in an AI system was influenced when it was recommended by another person.  
 11 Here, there was a stronger inclination towards reluctance among male participants  
 12 compared to female participants. Women participating in the study were more  
 13 inclined to consider others' trust as an important indicator for their own trust,  
 14 whereas men regarded their own trust in an AI system as more crucial.

15  
 16 **Figure 4.** *Diagram showing the single difference for non-safety-critical AI systems*

To use an AI system, the handling of the system must be easy and quick to learn.



18  
 19

1        Regarding questions related to safety-critical AI systems, no significant  
 2 differences between men and women were observed for any of the questions.  
 3 Similarly, for questions concerning non-safety-critical AI systems, only one  
 4 question yielded a significant difference. This question addressed the respectful  
 5 behavior of a chatbot towards its users. Here, female participants found it more  
 6 important than male participants for the chatbot to behave respectfully. A larger  
 7 proportion of male participants indicated that the chatbot's behavior towards the user  
 8 was less important, as more individuals stated they would still use the chatbot.  
 9 Figure 4 shows these observations. These slight discrepancies between the two  
 10 groups suggest that gender does not have huge implications on responses regarding  
 11 AI systems.

12

### 13 *Differences between Trust in Humans and in AI Systems*

14

15        At the beginning of the questionnaire, general questions about the trust of the  
 16 study participants were recorded. These questions related to the previously  
 17 mentioned questions about the social factors. They were derived from interpersonal  
 18 trust between two individuals. The questions about interpersonal trust as well as the  
 19 questions about social factors can accordingly be compared to determine a  
 20 difference between trust in humans and AI systems. As with the comparisons  
 21 between autonomous vehicles and chatbots previously, paired T-tests were  
 22 conducted for both appropriate questions to identify significant differences between  
 23 the responses. The following list shows the compared questions, as they differ  
 24 slightly in formulation, on the left side of the hyphen, the reference is to humans,  
 25 while on the right side, the reference is to AI systems. If the question is identical for  
 26 humans and AI systems, it was listed only once. This is the case for the pairs six,  
 27 eleven, and twelve.

28

29

- 30        1. whose actions I can comprehend. - whose algorithms I can comprehend.
- 31        2. who explains their behavior to me. - whose decision-making process I can  
 32 perceive.
- 33        3. whom I have known for a long time. - with whom I have interacted frequently.
- 34        4. who have acted reliably before. - which deliver reliable results.
- 35        5. who are responsible for their actions. - whose responsibilities I can assess.
- 36        6. who proceeds with great accuracy and care.
- 37        7. whom I perceive as competent. - which competently accomplish assigned  
 38 tasks.
- 39        8. who have previously executed a given task properly. - which have previously  
 40 delivered correct results.
- 41        9. who are like me. - which behave similarly to humans.
- 42        10. in whom I can empathize. - which are like AI systems I trust.
- 43        11. who have been recommended to me by other individuals.
- 44        12. whom I have trusted in the past.

45

1 Table 2 provides an overview of the paired T-tests. From this table, it is evident  
 2 that there are significant differences between trust in humans and AI systems. The  
 3 analysis revealed that there are only four questions where no significant differences  
 4 were found. It can be observed that no significant differences were found for pairs  
 5 two, four, ten, and twelve. Consequently, eight questions, concerning social factors,  
 6 differ in terms of trust between humans and AI systems. This indicates that a  
 7 different type of trust can be assumed when a person interacts with another person  
 8 than with an AI system.

9  
 10 **Table 2.** *T-test for comparing trust in humans and AI systems*

Pairs	T	One-Sided p	Two-Sided p
Pair 1	7.132	<.001	<.001
Pair 2	1.720	.043	.087
Pair 3	4.691	<.001	<.001
Pair 4	-2.808	.033	.005
Pair 5	-5.449	<.001	<.001
Pair 6	-5.966	<.001	<.001
Pair 7	-5.351	<.001	<.001
Pair 8	-8.565	<.001	<.001
Pair 9	10.310	<.001	<.001
Pair 10	.830	.204	.407
Pair 11	-5.382	<.001	<.001
Pair 12	1.171	.121	.243

11  
 12  
 13 Upon closer examination of the questions where no difference was detected, it  
 14 can be noted that understanding the behavior of both humans and AI systems is  
 15 similarly important. Likewise, past experiences in dealing with both humans and AI  
 16 systems hold a similar significance for the participants of this study. Also,  
 17 comparability and thus proximity to known and established trust structures  
 18 influence interactions with other humans or other AI systems. When interacting with  
 19 other humans, one's own person serves as a point of comparison, whereas with AI  
 20 systems, similar other AI systems are used. The last question reaffirms the  
 21 importance of experiences, as positive experiences in the past can have a positive  
 22 impact on trust in the future. All these questions reaffirm the previously mentioned  
 23 importance of trust. Previous experiences, similar systems, reliable results, and a  
 24 transparent decision-making process are key establishing trust in an AI system.

25 In summary, the high number of significant differences between humans and  
 26 AI systems suggests that trust must be viewed differently. Although the social  
 27 factors were developed with the idea of trust between two people, this comparison  
 28 illustrates well that while there are some similarities, a different trust relationship  
 29 must be assumed. However, the relationship between two people and the  
 30 experiences gained from it still provide clues to attempt an approximation of trust  
 31 between humans and AI systems.



## 1 Conclusion

2  
3 Based on the results of the study and the subsequent analysis, a closer  
4 examination of the individual factor groups should be the next step. Here, different  
5 types of AI systems are particularly important. It might be possible to establish  
6 hierarchies of factors based on the type of system. This would provide a good  
7 overview of the key factors for specific types of AI systems. Furthermore, another  
8 study with a more detailed examination of age groups could provide additional  
9 information regarding younger and older users.

10 Further interesting starting points can be found in the factors of transferability  
11 and behavior. Regarding transferability, it could be investigated whether human-like  
12 chatbots tend to increase trust, while human-like autonomous vehicles might  
13 decrease it. Additionally, a general direction of trust for the behavior of different AI  
14 systems to foster trust could be identified. In terms of behavior, aspects of privacy  
15 and ethics were examined. However, trust in an AI system would need to be  
16 assumed here to ensure that data misuse is avoided.

17 Furthermore, specific questions could be asked regarding the topic of  
18 Explainable AI. Especially the results from the comparison of human trust  
19 relationships to the trust relationship between a human and a given AI system  
20 yielded an emphasis on a transparent decision-making process. Here, respondents  
21 could first be provided with more detailed information and approaches through  
22 Explainable AI algorithms before questions are asked. This could help assess the  
23 effectiveness of specific methods. It would also be interesting to learn how  
24 explainability and transparency are desired by users. The study revealed that the  
25 option for reference was already sufficient for a large portion of the participants,  
26 indicating that an information overload may not be desired. This would mean that  
27 only in cases of uncertainty the results of an algorithm would be verified to uncover  
28 potential problems or errors in the algorithm's solution process.

29 The analysis of trust between humans and AI systems generally offers great  
30 potential and allows for a better assessment of which aspects of a system should be  
31 given special attention in the development process. This is especially important for  
32 technologies intended for commercial use. Different people have different notions  
33 of trust, and therefore, consideration should be given to how trust in a new  
34 technology can be established or improved. As mentioned at the beginning of this  
35 paper, it cannot be assumed that all users fully understand the technology. Trust is  
36 needed to fill this "knowledge gap", enabling the use of technology without  
37 concerns.

## 40 References

- 41  
42 Alarcon GM, Gibson AM, Jessup SA, Capiola A, Raad H, Lee MA (2020) Effects of  
43 Reputation, Organization, and Readability on Trustworthiness Perception of Computer  
44 Code. *Human-Computer Interaction: Human Values and Quality of Life*: 367-825.  
45 Bach TA, Khan A, Hallock H, Beltrão G, Sousa S (2022) A Systematic Literature Review  
46 of User trust in AI-Enabled Systems: An HCI Perspective. *International Journal of*  
47 *Human-Computer Interaction* 40(5): 1251-1266.

- 1 Bonezzi A, Ostinelli M, Melzner J (2022) The Human Black-Box: The Illusion of  
 2 Understanding Human Better Than Algorithmic Decision-Making. *Journal of*  
 3 *Experimental Psychology General* 151(9): 2250-2258.
- 4 Carleton RN (2016) Into the unknown: A review and synthesis of contemporary models  
 5 involving  
 6 uncertainty. *Journal of Anxiety Disorders*: 30-43.
- 7 Devitt SK (2018) Trustworthiness of Autonomous Systems, Foundation of Trusted  
 8 Autonomy. *Foundations of Trusted Autonomy*: 161-184
- 9 European Commission: Directorate-General for Communications Networks, Content and  
 10 Technology (2019) *Ethics guidelines for trustworthy AI*. Publications Office.
- 11 Ferronato P, Bashir M (2020) An Examination of Dispositional Trust in Human and  
 12 Autonomous System Interactions. *Human-Computer Interaction: Human Values and*  
 13 *Quality of Life, Thematic Area, HCI 2020: Proceedings Part III*: 420-435
- 14 Floridi L, Cowls J, Beltrametti M, Chatila R, Chazerand P, Dignum V, Luetge C, Madelin  
 15 R, Pagallo U, Rossi F, Schafer B, Valcke P, Vayena E (2019) AI4People – An Ethical  
 16 Framework for a Good AI Society: Opportunities, Risks, Principles, and  
 17 Recommendations. *Minds and Machines* 28: 689-707.
- 18 Goodman, B. & Flaxman S. (2016). EU regulations on algorithmic decision making and a  
 19 “right to explanation”: *AI Magazine* 38(3): 50-57.
- 20 Hancock PA, Billings DR, Schaefer KE (2011) Can You Trust Your Robot?. *Ergonomics in*  
 21 *Design The Quarterly of Human Factors Applications* 19(3): 24-29.
- 22 Hoff, K. A. & Bashir, M. (2015) Trust in Automation: Integrating Empirical Evidence on  
 23 Factors that Influence Trust. *Human Factors: The Journal of the human Factors and*  
 24 *Ergonomics Society* 57(3): 403-434.
- 25 Ikkatai Y, Hartwig T, Takanashi N, Yokoyama HM (2022) Octagon Measurement: Public  
 26 Attitudes toward AI Ethics. *International Journal of Human-Computer Interactions*  
 27 38(17): 1589-1606.
- 28 Jacovi A, Marasovic A, Miller T, Goldberg Y (2021) Formalizing Trust in Artificial  
 29 Intelligence: Prerequisites, Causes and Goals of Human Trust in AI. *In Proceedings of*  
 30 *the 2021 ACM Conference on Fairness, Accountability, and Transparency*: 624–635.
- 31 Kaplan AD, Kessler TT, Brill JC, Hancock PA (2021) trust in Artificial Intelligence: Meta-  
 32 Analytic Findings. *Human Factors: The Journal of the Human Factors and Ergonomic*  
 33 *Society* 65(2).
- 34 Kusche I (2024) Possible harms of artificial intelligence and the EU AI act: fundamental  
 35 rights and risk. *Journal of Risk Research*: 1-14.
- 36 Luhmann N (1968) Trust. A Mechanism for the Reduction of Social Complexity. (5th ed).  
 37 UVK Wiesbaden.
- 38 OECD (2017) OECD Guidelines on Measuring Trust, *OECD Publishing*, Paris.
- 39 Parasuraman R, Riley V (1997) Human Use of Automation: Use, Misuse, Disuse, Abuse.  
 40 *Human Factors: The Journal of the Human Factors and Ergonomics Society* 39(2):  
 41 230-253.
- 42 Smithson M (2018) Trusted Autonomy Under Certainty. *Foundation of Trusted Autonomy:*  
 43 *Studies in Systems, Decision and Control* 117: 185-201.